

Résumé du Cours de Statistique Descriptive

Yves Tillé 18 janvier 2008

Objectif et moyens

Objectifs du cours – Apprendre les principales techniques de statistique descriptive univariée et bivariée.

- Être capable de mettre en oeuvre ces techniques de manière appropriée dans un contexte donné.
- Être capable d'utiliser les commandes de base du langage R. Pouvoir appliquer les techniques de statistiques descriptives au moyen du langage R.
- Références Dodge Y. (2003), *Premiers pas en statistique*, Springer.

Droesbeke J.-J. (1997), *Éléments de statistique*, Editions de l'Université libre de Bruxelles/Ellipses.

Moyens – 2 heures de cours par semaine. – 2 heures de TP par semaine, réparties en TP théoriques et appliqués.

Le langage R – Shareware : gratuit et installé en 10 minutes. – Open source (on sait ce qui est réellement

Table des matières

1 Variables, données statistiques, tableaux, effectifs	5
1.1 Définitions fondamentales	5
1.1.1 La science statistique	5
1.1.2 Mesure et variable	5
1.1.3 Typologie des variables	5
1.1.4 Série statistique	6
1.2 Variable qualitative nominale	6
1.2.1 Effectifs, fréquences et tableau statistique	6
1.2.2 Diagramme en secteur et diagramme en barres	7
1.3 Variable qualitative ordinale	7
1.3.1 Le tableau statistique	7
1.3.2 Diagramme en secteurs	9
1.3.3 Diagramme en barres des effectifs	10
1.3.4 Diagramme en barres des effectifs cumulés	10
1.4 Variable quantitative discrète	11
1.4.1 Le tableau statistique	11
1.4.2 Diagramme en bâtonnets des effectifs	12
1.4.3 Fonction de répartition	12
1.5 Variable quantitative continue	12
1.5.1 Le tableau statistique	12
1.5.2 L'histogramme des effectifs	14
1.5.3 La fonction de répartition	15

2 Statistique descriptive univariée 17

2.1 Paramètres de position	17
2.1.1 Le mode	17
2.1.2 La moyenne	17
2.1.3 Remarques sur le signe de sommation	
L'étendue	24
2.2 La distance interquartile	24
2.4.2 Coefficient d'asymétrie de Yule	27
2.4.3 Coefficient d'asymétrie	27

2.5 Paramètre d'aplatissement (kurtosis).	28
2.6 Changement d'origine et d'unité.	29
2.7 Moyennes et variances dans des groupes.	29
2.8 Diagramme en tiges et feuilles.	31
2.9 La boîte à moustaches.	31
3 Statistique descriptive bivariée	35
3.1 Séries statistiques bivariées.	35
3.2 Deux variables quantitatives.	35
3.2.1 Représentation graphique de deux variables.	35
3.2.2 Analyse des variables.	36
3.2.3 Covariance.	36
3.2.4 Corrélation.	37
3.2.5 Droite de régression.	37
3.2.6 Résidus et valeurs ajustées.	41
3.2.7 Sommes de carrés et variances.	41
3.2.8 Décomposition de la variance.	42
3.3 Deux variables qualitatives.	43
3.3.1 Données observées.	43
3.3.2 Tableau de contingence.	44
3.3.3 Tableau des fréquences.	44
3.3.4 Profils lignes et profils colonnes.	45
3.3.5 Effectifs théoriques et khi-carré.	46
4 Théorie des indices, mesures d'inégalité	51
4.1 Nombres indices.	51
4.2 Définition.	51
4.2.1 Propriétés des indices.	52
4.2.2 Indices synthétiques.	52
4.2.3 Indice de Laspeyres.	52
4.2.4 Indice de Paasche.	53
4.2.5 L'indice de Fisher.	53
4.2.6 L'indice de Sidgwick.	54
4.2.7 Indices chaînés.	54
4.3 Mesures de l'inégalité.	54
4.3.1 Introduction.	54
4.3.2 Courbe de Lorenz.	55
4.3.3 Indice de Gini.	56
4.3.4 Indice de Hoover.	56
4.3.5 Quintile et Decile share ratio.	56
Exemples.	59
5.2 Description de la tendance.	59
6.3.2 Opérateur différence.	65

5.3.3 Différence saisonnière	67
5.4 Filtres linéaires et moyennes mobiles	69
5.4.1 Filtres linéaires.	69
5.4.2 Moyennes mobiles : définition.	70
5.4.3 Moyenne mobile et composante saisonnière	70
5.5 Moyennes mobiles particulières	71
5.5.1 Moyenne mobile de Van Hann.. . . .	71
5.5.2 Moyenne mobile de Spencer.	71
5.5.3 Moyenne mobile de Henderson	71
5.5.4 Médianes mobiles.. . . .	72
5.6 Désaisonnalisation	72
5.6.1 Méthode additive	72
5.6.2 Méthode multiplicative.. . . .	73
5.7 Lissage exponentiel	73
5.7.1 Lissage exponentiel simple	73
5.7.2 Lissage exponentiel double	76
6 Calcul des probabilités et variables aléatoires	83
6.1 Probabilités	83
6.1.1 Événement	83
6.1.2 Opérations sur les événements	83
6.1.3 Relations entre les événements	84
6.1.4 Ensembles des parties d'un ensemble et système complet.	84
6.1.5 Axiomatique des Probabilités	84
6.1.6 Probabilités conditionnelles et indépendance	87
6.1.7 Théorèmes des probabilités totales et théorème de Bayes	87
6.2 Analyse combinatoire	88
6.2.1 Introduction	88
6.2.2 Permutations (sans répétition)	88
6.2.3 Permutations avec répétition	89
6.2.4 Arrangements (sans répétition)	89
6.2.5 Combinaisons	89
6.3 Variables aléatoires	90
6.3.1 Définition	90
6.4 Variables aléatoires discrètes.	90
6.4.1 Définition, espérance et variance	90
6.4.2 Variable indicatrice de Bernoulli	91
6.4.3 Variable binomiale	91
6.4.4 Variable de Poisson	96
6.5 Variables aléatoires continues	96
6.5.1 Cas continu	96
6.5.2 Indépendance	96
7 Tables statistiques	1074

Chapitre 1

Variables, données statistiques, tableaux, effectifs

1.1 Définitions fondamentales 1.1.1 La science statistique – Méthodes scientifiques

– Etymologiquement: science de l'état. – La statistique s'applique à la plupart des disciplines : agronomie, biologie, démographie, économie, sociologie, linguistique, psychologie, ...

1.1.2 Mesure et variable – On s'intéresse à des *unités statistiques*

des entreprises, des ménages. En sciences humaines, on s'intéresse dans la plupart des cas à un nombre fini d'unités. – Sur ces unités, on mesure un caractère ou une *variable*, le chiffre d'affaires de l'entreprise, le revenu du ménage, l'âge de la personne, la catégorie socio-professionnelle d'une personne. On suppose que la variable prend toujours une seule valeur sur chaque unité. Les variables sont désignées par simplicité par une lettre (X, Y, Z). – Les *valeurs possibles* de la variable, sont appelées *modalités*. – L'ensemble des valeurs possibles fait de pouvoir ou non ordonner les modalités est parfois discutable. Par exemple : dans les catégories sociales. Une variable est dite continue, si l'ensemble des valeurs possibles est continu. **Remarque 1.1** Ces définitions

Exemple 1.1 Les modalités de la variable sexe sont *masculin* (M) et *f minin* (F). Le domaine de la variable est $\{M, F\}$.

Exemple 1.2 Les modalités de la variable nombre d'enfants par famille sont 0,1,2,3,4,5, une variable quantitative discr te.

1.1.4 S rie statistique

On appelle *s rie statistique* la suite des valeurs prises par une variable X sur n observations.

Le nombre d'unit es d'observation est not 

Les valeurs de la variable X sont not es

$$x_1, \dots, x_i, \dots, x_n.$$

Exemple 1.3 On s'int resse   la variable " tat-civil". On a la s rie statistique des valeurs prises par X sur 20 personnes. La codification est C: c libataire, M : mari e(e), V : veuf(ve), D : divorc e. Le domaine de

=====

=====

=====

$M M D C C M C C C M C M V M V D C C C M$ Ici, $n=20$, $x_1 = M, x_2 = M, x_3 = D, x_4 = C, x_5 = C, \dots, x_{20} = M$. **1.2 Var**
l'effectif de la modalit  x_j . La fr quence d'une modalit  est l'effectif divis  par le nombre d'unit es d'obs

—

=====

=====

=====

=====

En langage R

```
>X=c('Mari e(e)','Mari e(e)','Divorc e(e)','C elibataire','C elibataire','Mari e(e)','C elibataire',  
      'C elibataire','C elibataire','Mari e(e)','C elibataire','Mari e(e)','Veuf(ve)','Mari e(e)',  
      'Veuf(ve)','Divorc e(e)','C elibataire','C elibataire','C elibataire','Mari e(e)')  
>T1=table(X)> V1=c(T1)> data.frame(Eff=V1,Freq=V1/sum(V1))
```

Eff Freq C elibataire 9 0.45 Divorc e(e) 2 0.10 Mari e(e) 7 0.35 Veuf(ve) 2 0.10

1.2.2 Diagramme en secteurs

Le tableau statistique peut  tre repr sent e par un diagramme en barres ou en secteurs (ou camembert ou *pie chart* en anglais) (voir Figures 1.1 et 1.2). C elibataire

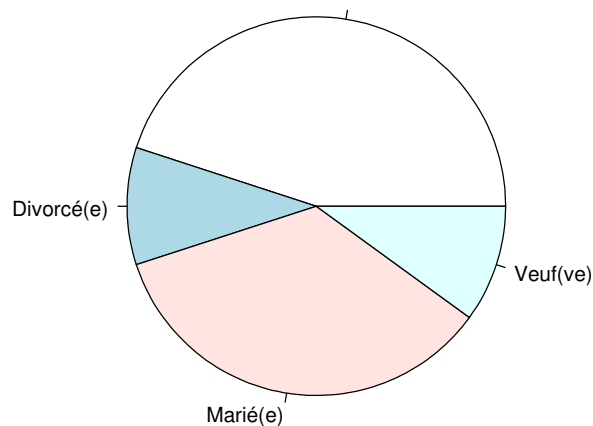


Fig. 1.1 - Diagramme en secteurs

En langage R > pie(T1,radius=1.0)

En langage R > barplot(T1)

1.3 Valeurs



Fig.1.2-Diagramme enbarres

La notation $x_1 < x_2$ se lit x_1 précède x_2 . Si la variable est ordinale, on peut calculer les effectifs cumulés:

$$N_j = \sum_{k=1}^j n_k.$$

On a $N_1 = n_1$ et $N_J = n$. On peut également calculer les fréquences cumulées

$F_j = N_j / n = \sum_{k=1}^j f_k$. **Exemple 1.5** On interroge 50 personnes sur leur dernier diplôme obtenu (variable Y).

La codification a été faite selon le Tableau 1.1. On a obtenu la série statistique présentée dans le tableau 1.2. Finalement, on obtient le tableau statistique complet présenté dans le Tableau 1.3. Tab.1.1-Codification de la variable Y Dernier diplôme obtenu

En langage R > YY=c("Sd","Sd","Sd","Sd","P","P","P","P","P","P","P","P","P","P","Se","Se","Se","Se","Se",

Tab. 1.3- Tableau statistique complet

x_j	n_j	N_j	f_j	F_j
Sd	4	4	0.08	0.08
P	11	15	0.22	0.30
Se	14	29	0.28	0.58
Su	9	38	0.18	0.76
U	12	50	0.24	1.00
	50		1.00	

```
T2=table(YF)V2=c(T2)> data.frame(Eff=V2,EffCum=cumsum(V2),Freq=V2/sum(V2),FreqCum=cumsum(V2),
Eff EffCum Freq FreqCumSd 4 4 0.08 0.08P 11 15 0.22 0.30Se 14 29 0.28 0.58Su 9 38 0.18 0.76U 12 50 0.24 1.00
```

Figure 1.3). Sd

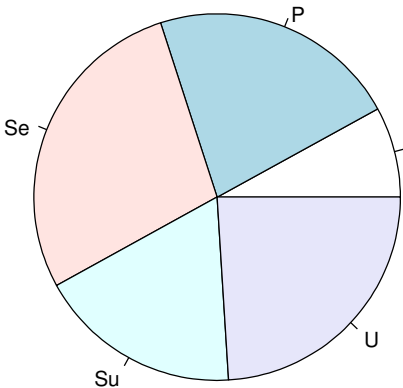


Fig.1.3- Diagrammeen secteursdesfréquences**En langage R**> pie(T2,radius=1)9

1.3.3 Diagramme en barres des effectifs

Les effectifs d'une variable qualitative sont représentés au moyen d'un diagramme en barres (voir Figure 1.4). *sd*

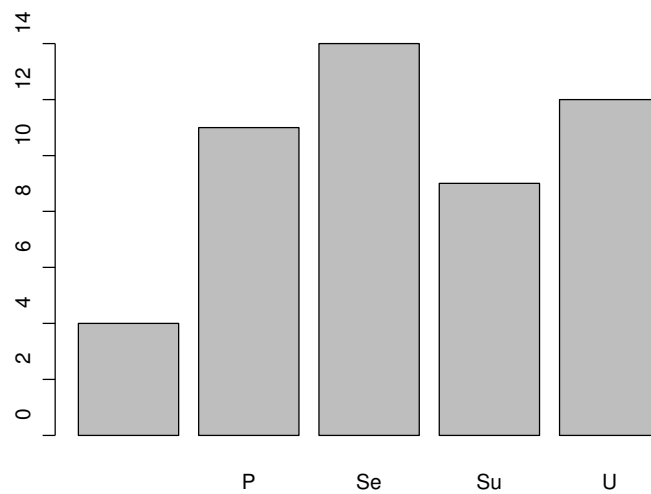


Fig. 1.4- Diagramme en barres des effectifs

En langage R `> barplot(T2)` **1.3.4 Diagramme en barres des effectifs cumulés** *sd*

Figure 1.5). *sd*

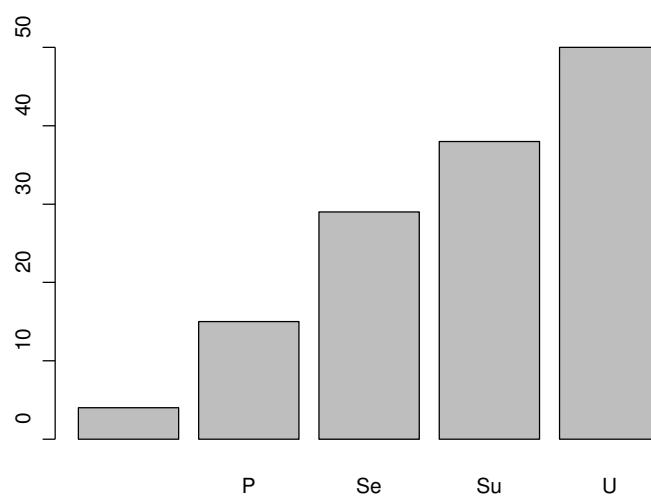


Fig. 1.5-Diagramme en barres des effectifs cumulés

En langageR

```
>T3=cumsum(T2)
>barplot(T3)
```

1.4 Variable quantitative discontinue

1.4.1 Le tableaustatistique

Une variable discrète a un domaine dénombrable.

Exemple 1.6 Un quartier est compos e de 50 m ges, et la variable Z repr sente le nombre de personnes par m n ge. Les valeurs de la variable sont 1 1 1 1 1

	2	2	2	2	2
2 22 233		3	3	3	3
3 33 33 33			3	3	4
4 4 4 4 4 444				4	5
5 5 5 5 5 66 68					8

Comme pour les variables qualitatives ordinales, on peut calculer les effectifs, les effectifs cumulés, les fréquences, les fréquences cumulées. À nouveau, on peut construire le tableau statistique:

$$x_j \ n_j \ N_j \ f_j \ F_j | 5 \ 5 \ 0.10 \ 0.102 \ 9 \ 14 \ 0.18 \ 0.283 \ 15 \ 29 \ 0.30 \ 0.584 \ 10 \ 39 \ 0.20 \ 0.785 \ 6 \ 45 \ 0.12 \ 0.906 \ 3 \ 48 \ 0.06 \ 0.192$$

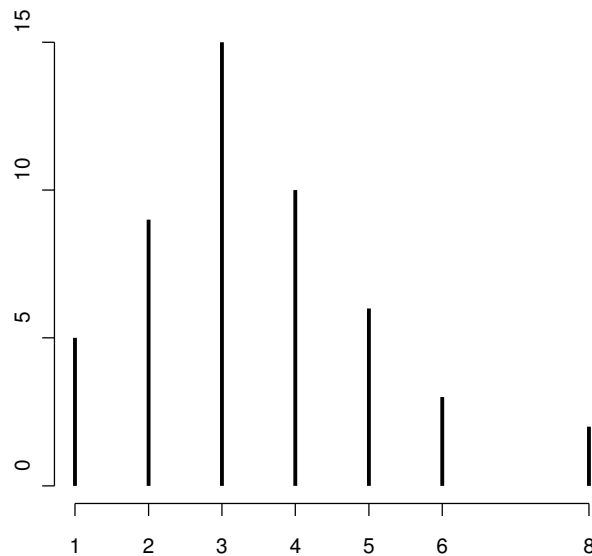


Fig. 1.6 – Diagramme en bâtonnets des effectifs pour une variable quantitative discrète

1.4.2 Diagramme en bâtonnets des effectifs

Quand la variable est discrète, les effectifs sont

En langage R `> plot(T4,type="h",xlab="",ylab="",main="",frame=0,lwd=3)`

1.4.3 Fonction de répartition

est définie

en Figure 1.7, est définie de R dans $[0, 1]$ et vaut : $F(x) = \begin{cases} 0 & x < x_1 \\ F_j & x_j \leq x < x_{j+1} \\ 1 & x_j \leq x \end{cases}$. **En langage R** `> plot(ecdf(T4))`

discrètes. Cependant, il est souvent intéressant de procéder à des regroupements en classes pour faire



tediscr`

52 153 153154 154 154 155 1551
distribution group'ee. On note, de

- N_j l'effectif cumulé de la classe j ,
- f_j la fréquence de la classe j ,
- F_j la fréquence cumulée de la classe j .

En langage R > S=c(152,152,152,153,153,154,154,154,155,155,156,156,156,156,

+157,157,157,158,158,159,159,160,160,160,161,160,160,161,162,

+162,162,163,164,164,164,164,165,166,167,168,168,168,169,169,

+ 170,171,171,171,171)> T5=table(cut(S, breaks=c(151,155,159,163,167,171)))

>T5c=c(T5)> data.frame(Eff=T5c,EffCum=cumsum(T5c),Freq=T5c/sum(T5c),FreqCum=cumsum(T5c/sum(T5c)))

Eff EffCum Freq FreqCum(151,155] 10 10 0.20 0.20(155,159] 12 22 0.24 0.44(159,163] 11 33 0.22 0.66(163,167] 6 39 0.12 0.78(167,171] 9 48 0.18 1.00

hauteur) représente l'effectif. La hauteur h_j du rectangle correspondant à la classe j est donc donné

$$h_j = n_j c_{j+1} - c_j . 0$$

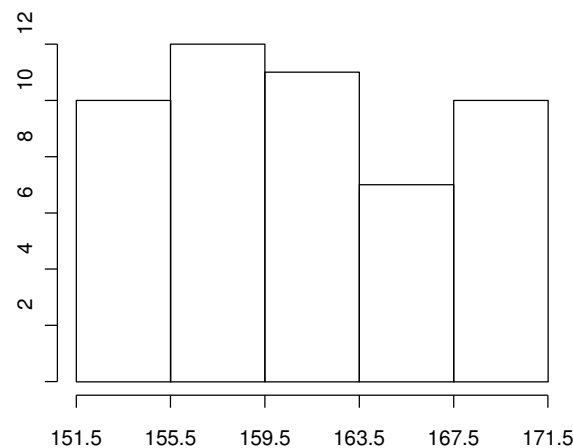


Fig.1.8-Histogrammedeseffectifs**En langage R** > hist(S,breaks=c(151.5,155.5,159.5,163.5,167.5,171.5),xlab="")

Si les deux dernières classes sont égales, comme dans la Figure 1.9, la surface du dernier rectangle est égale à la surface des deux premiers rectangles de l'histogramme de la Figure 1.8.

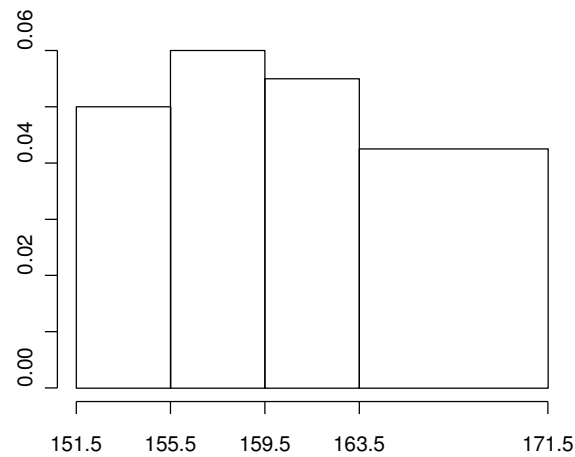
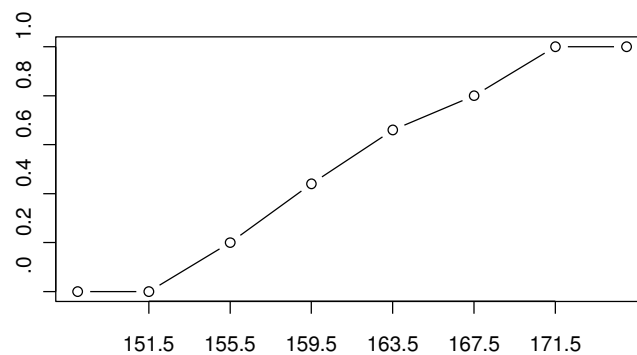


Fig. 1.9- Histogramme des effectifs avec les deux dernières classes égales

En langage R `> hist(S,breaks=c(151.5,155.5,159.5,163.5,171.5),xlab="",ylab="",main="",xaxt = "n")> ax`



En langageR

```
>y=c(0,0,cumsum(T5c/sum(T5c)),1)
>x=c(148,151.5,155.5,159.5,163.5,167.5,171.5,175)
> plot(x,y,type="b",xlab="",ylab="",xaxt ="n")
> axis(1, c(151.5,155.5,159.5,163.5,167.5,171.5))
```

Chapitre 2

Statistique descriptive univariée

2.1 Paramètres de position

2.1.1 Le mode

Le mode est la valeur distincte correspondant à l'effectif le plus élevé ; il est noté x_M .

Si on reprend la variable 'Etat civil', dont le tableau statistique est le suivant :

x_j n_j f_j C 9 0.45 M 7 0.35 V 2 0.10 D 2 0.10 $n = 20$ 1 le mode est C : célibataire.

Remarque 2.1 – Le mode peut être

pondant à l'effectif le plus élevé).

2.1.2 La moyenne

La *moyenne* ne peut être définie que sur une

—

On peut aussi faire les calculs avec les valeurs distinctes et les effectifs se résume à :

x_j	n_j
0	2
1	3
2	1
3	1
4	1
	8

$$\begin{aligned}\bar{x} &= \frac{2 \times 0 + 3 \times 1 + 1 \times 2 + 1 \times 3 + 1 \times 4}{8} \\ &= \frac{3 + 2 + 3 + 4}{8}\end{aligned}$$

= 1.5. **Remarque 2.2** La moyenne n'est pas nécessairement une valeur possible.

En langage R $E=c(0,0,1,1,1,2,3,4)$ $n=length(E)$ $\bar{x}=sum(E)/n$ $\bar{x}=mean(E)$

2.1.3 Remarques sur la

$x_{31} + x_{32} \ (i = 3)$

5. On peut exclure une valeur de l'indice.

$$\sum_{\substack{j=1 \\ j \neq 3}}^5 x_j = x_1 + x_2 + x_4 + x_5.$$

Propriété 2.11. Somme d'une constante $\sum_{i=1}^n a = a + a + \dots + a$

$$\sum_{i=1}^n a = na \quad (a \text{ constante}).$$

Exemple $\sum_{i=1}^5 3 = 3 + 3 + 3 + 3 + 3 = 5 \times 3 = 15.$

2. Mise en évidence $\sum_{i=1}^n ax_i = a \sum_{i=1}^n x_i$ (a constante).

Exemple $\sum_{i=1}^3 2 \times i = 2(1 + 2 + 3) = 2 \times 6 = 12.$ 3. Somme des n premiers entiers $\sum_{i=1}^n i = 1 + 2 + 3 + \dots + n = n(n+1)/2$

2.1.4 Moyenne géométrique

Si $x_i \geq 0$, on appelle moyenne géométrique la quantité

$$G = \sqrt[n]{x_1 \times x_2 \times \dots \times x_n} = \exp \frac{1}{n} \sum_{i=1}^n \log x_i.$$

La moyenne géométrique s'utilise, par exemple, quand on veut calculer la moyenne des taux d'intérêt.

Exemple 2.3 Supposons que les taux d'intérêt pour 4 années consécutives soient respectivement de 5, 10, 15, et 10%. Que va-t-on obtenir après 4 ans si je place 100 francs?

- Après 1 an on a, $100 \times 1.05 = 105$ Fr.
- Après 2 ans on a, $100 \times 1.05 \times 1.1 = 115.5$ Fr.
- Après 3 ans on a, $100 \times 1.05 \times 1.1 \times 1.15 = 132.825$ Fr.
- Après 4 ans on a, $100 \times 1.05 \times 1.1 \times 1.15 \times 1.1 = 146.1075$ Fr.

Si on calcule la moyenne arithmétique des taux on obtient

$$\bar{x} = \frac{1.05 + 1.10 + 1.15 + 1.10}{4} = 1.10.$$

Si on calcule la moyenne géométrique des taux, on obtient $G = (1.05 \times 1.10 \times 1.15 \times 1.10)^{1/4}$

$$= 1.099431377.$$

Le bon taux moyen est bien G et non \bar{x} , car si on applique 4 fois le taux moyen G aux 100 francs, on obtient

$$100 \text{ Fr} \times G^4 = 100 \times 1.099431377^4 = 146.1075 \text{ Fr.}$$

2.1.5 Moyenne harmonique Si $x_i \geq 0$, on appelle moyenne harmonique la quantité $H = n / \sum_{i=1}^n 1/x_i$. Il est facile de vérifier que si x_1, x_2, \dots, x_n sont des vitesses, la vitesse moyenne est donc $Moy = \frac{400}{20.8333} = 19.2 \text{ km/h.}$ - Si on calcule la moyenne

Remarque 2.3 Il est possible de montrer que la moyenne harmonique est toujours inférieure à la moyenne géométrique qui est toujours inférieure à la moyenne arithmétique.

$$H \leq G \leq \bar{x}.$$

2.1.6 Moyenne pondérée

Dans certains cas, on n'accorde pas le même poids à toutes les observations. Par exemple, si on calcule la moyenne des notes pour un programmeur, on peut pondérer les notes de l'étudiant par le nombre de crédits ou par le nombre d'heures de chaque cours. Si $w_i = 1, \dots, n$ sont les poids associés à chaque observation, alors la moyenne pondérée se définit par :

$$\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}.$$

Exemple 2.5 Supposons que les notes soient pondérées par le nombre de crédits, et que les notes de l'étudiant soient les suivantes : Note 5 4 3 6 5 Crédits 6 3 4 3 4 La moyenne pondérée des notes par les

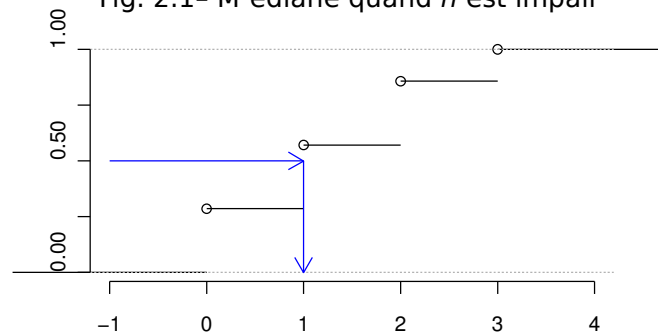
$$\frac{5 \cdot 6 + 4 \cdot 3 + 3 \cdot 4 + 6 \cdot 3 + 5 \cdot 4}{6 + 3 + 4 + 3 + 4} = \frac{92}{20} = 4.6.$$

2.1.7 La médiane

La médiane, notée $x_{1/2}$, est une valeur centrale de la série statistique obtenue d'une

- On trie la série statistique par ordre croissant des valeurs observées. Avec la série observée : 3 2 1 4 5, la médiane est la valeur qui se trouve à la position 1/2 de la répartition pour la valeur 1/2 : $x_{1/2} = F_{-1}(0.5)$. **En langage R** : `median(x)`

Fig. 2.1- M'ediane quand n est impair



`x=c(0 , 0 , 1 , 1 , 2 , 2 , 3)median(x)plot(ecdf(x),xlab="",ylab="",main="",frame=FALSE,yaxt = "n")`

`axis(2, c(0.0,0.25,0.50,0.75,1.00))arrows(-1,0.5,1,0.50,length=0.14,col="blue")`

`arrows(1,0.50,1,0,length=0.14,col="blue")` - Si n est pair, deux valeurs se trouvent au milieu de la série (ici avec $n=8$)

0 0 1 1 2 2 3 4

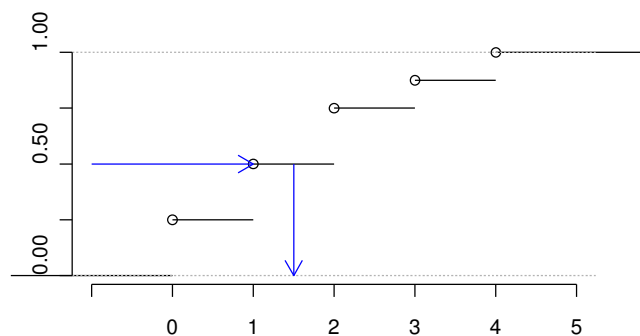
↑ ↑ La médiane est alors la moyenne de ces deux valeurs : $x_{1/2} = 1 + 2 = 1.5$. La Figure 2.2 montre la

—

médiane peut toujours

être définie comme l'inverse de la fonction de répartition pour la valeur $1/2$: $x_{1/2} = F^{-1}(0.5)$. Cependant

exactement à un 'palier'. Fig. 2.2- Médiane quand n est pair-1



Exemple en langage R `x=c(0 , 0 , 1 , 1 , 2 , 2 , 3 , 4)median(x)plot(ecdf(x),xlab="",ylab="",main="",frame=FALSE,yaxt = "n")`

En général on note

$$X_{(1)}, \dots, X_{(i)}, \dots, X_{(n)}$$

la série ordonnée par ordre croissant. On appelle cette ordonnée la statistique d'ordre. Cette notation, très usuelle en statistique, permet d'exprimer la médiane de manière très synthétique.

- Si n est impair

$$X_{1/2} = X_{(\frac{n+1}{2})}$$

- Si n est pair

$$X_{1/2} = \frac{1}{2} X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)}$$

Remarque 2.4 La médiane peut être calculée sur des variables quantitatives et sur des variables qualitatives ordinales.

2.1.8 Quantiles

La notion de quantile d'ordre p (où $0 < p < 1$) généralise la médiane. Formellement un quantile est donné

par l'inverse de la fonction de répartition : $x_p = F^{-1}(p)$.

Si la fonction de répartition était continue et strictement croissante, la définition du quantile serait sans équivoque. La fonction de répartition est cependant discontinue par palier. Quand la fonction de répartition est par palier, il existe au moins 9 manières différentes de définir des quantiles selon que l'on fasse ou non une interpolation de la fonction de répartition. Nous présentons deux des méthodes les plus utilisées. Il ne faut pas s'étonner de voir les valeurs des quantiles différer légèrement d'un logiciel statistique à l'autre.

- Si np est un nombre entier, alors $x_p = \frac{1}{2} X_{(np)} + \frac{1}{2} X_{(np+1)}$. - Si np n'est pas un nombre entier, alors $x_p = X_{(np)}$.

Exemple : la statistique 12, 13, 15, 16, 18, 19, 22, 24, 25, 27, 28, 34 contenant 12 observations ($n = 12$). - Le premier quantile

- Le troisième quartile : Comme $np = 0.75 \times 12 = 9$ est un nombre entier, on a

$$x_{3/4} = \frac{x_{(9)} + x_{(10)}}{2} = \frac{25 + 27}{2} = 26.$$

En langage R $x=c(12,13,15,16,18,19,22,24,25,27,28,34)$

`quantile(x,type=2)` **Exemple 2.7** Soit la série statistique 12,13, 15, 16,18, 19, 22, 24,25, 27 contenant 10 observations.

- Le premier quartile : Comme $np = 0.25 \times 10 = 2.5$ n'est pas un nombre entier, on a

$$x_{1/4} = x_{(2.5)} = x_{(3)} = 15.$$

- La médiane : Comme $np = 0.5 \times 10 = 5$ est un nombre entier, on a

$$x_{1/2} = \frac{x_{(5)} + x_{(6)}}{2} = (18 + 19)/2 = 18.5.$$

- Le troisième quartile : Comme $np = 0.75 \times 10 = 7.5$ n'est pas un nombre entier, on a

$x_{3/4} = x_{(7.5)} = x_{(8)} = 24.$ **En langage R** $x=c(12,13,15,16,18,19,22,24,25,27)$ `quantile(x,type=2)` **2.2 Paramètres de position**

2.2.3 La variance

La *variance* est la somme des carrés des écarts à la moyenne divisée par le nombre d'observations:

$$s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Théorème 2.1 La variance peut aussi s'écrire

$$s_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2. \quad (2.1)$$

Démonstration $s_{2x} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 =$

$$\begin{aligned} &= \frac{1}{n} \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x} \frac{1}{n} \sum_{i=1}^n x_i + \bar{x}^2 \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}\bar{x} + \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \end{aligned}$$

La variance peut également être définie à partir des effectifs et des valeurs distinctes: *

$s_{2x} = \frac{1}{n} \sum_{j=1}^J n_j x_j^2 - \bar{x}^2$. La variance peut aussi s'écrire $s_{2x} = \frac{1}{n} \sum_{j=1}^J n_j x_j^2 - \bar{x}^2$. Quand on veut estimer une va

sélectionnée au hasard) de taille n , on utilise la variance "corrigée" divisée par $n-1$. $s_{2x} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
 $2+3+4+4+5+6+7+98 = 5,25$

$$s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$= \frac{1}{18} (2^2 + 3^2 + 4^2 + 4^2 + 5^2 + 6^2 + 7^2 + 9^2) - 5^2$$

$$= \frac{1}{18} (9 + 4 + 1 + 1 + 0 + 1 + 4 + 16)$$

= 368 = 4.5. On peut également utiliser la formule (2.1) de la variance,

ce qui nécessite moins de calculs (surtout quand la moyenne n'est pas un nombre entier).

$$s_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

$$= \frac{1}{18} (2^2 + 3^2 + 4^2 + 4^2 + 5^2 + 6^2 + 7^2 + 9^2) - 5^2$$

$$= \frac{1}{18} (4 + 9 + 16 + 16 + 25 + 36 + 49 + 81) - 25$$

= 2368 / 18 - 25 = 29.5 - 25 = 4.5. **En langage R** > x=c(2,3,4,4,5,6,7,9) > n=length(x) > s2=sum((x-mean(x))^2)/n

2.2.5 L'écart moyen absolu

L'écart moyen absolu est la somme des valeurs absolues des écarts à la moyenne divisée par le nombre d'observations :

$$e_{moy} = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|.$$

2.2.6 L'écart médian absolu

L'écart médian absolu est la somme des valeurs absolues des écarts à la médiane divisée par le nombre d'observations :

$$e_{med} = \frac{1}{n} \sum_{i=1}^n |x_i - x_{1/2}|.$$

2.3 Moments

Définition 2.2 On appelle moment à l'origine d'ordre $r \in \mathbb{N}$ le paramètre

$$m_r = \frac{1}{n} \sum_{i=1}^n x_i^r.$$

Définition 2.3 On appelle moment centré d'ordre $r \in \mathbb{N}$ le paramètre

$m_r = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^r$. Les moments généralisent la plupart des paramètres. On a en particulier - $m_1 = \bar{x}$

et l'aplatissement. **2.4 Paramètres de forme**

2.4.1 Coefficient d'asymétrie de Fisher normalisé par la distance interquartile : $A_F = \frac{x_{3/4} + x_{1/4} - 2x_{1/2}}{x_{3/4} - x_{1/4}}$. 27

2.4.3 Coefficient d'asymétrie de Pearson

Le coefficient d'asymétrie de Pearson est basé sur une comparaison de la moyenne et du moment standardisé par rapport à la médiane :

$$A_P = \frac{\bar{x} - M}{s_x}$$

Tous les coefficients d'asymétrie ont les mêmes propriétés, ils sont nuls si la distribution est symétrique, négatifs si la distribution est allongée à gauche (left asymmetry), et positifs si la distribution est allongée à droite (right asymmetry) comme montré dans la Figure 2.3.

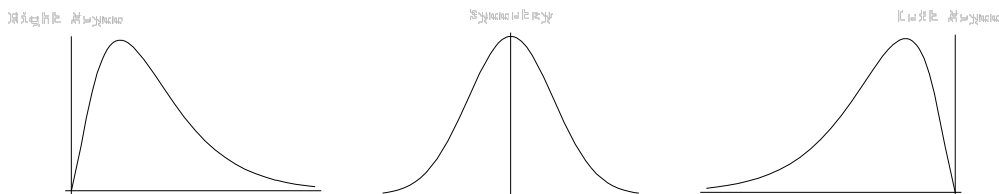


Fig. 2.3-Asymétrie d'une distribution

Remarque 2.6 Certaines variables sont toujours très asymétriques à droite, comme les revenus, les tailles des entreprises, ou des communes. Une méthode simple pour rendre une variable symétrique consiste alors à prendre le logarithme de cette variable.

2.5 Paramètre d'aplatissement (kurtosis)

Dans la Figure 2.4, on présente un exemple de deux distributions de même moyenne et de même

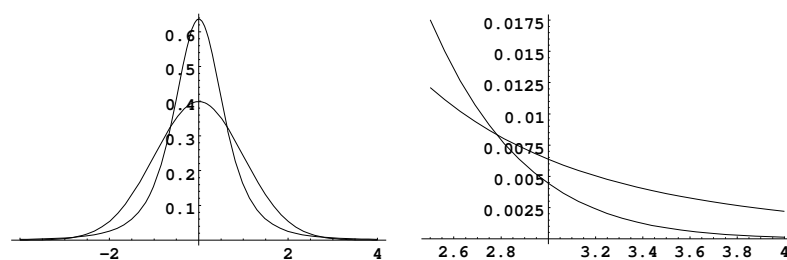


Fig. 2.4- Distributions mésokurtique et leptokurtique

2.6 Changement d'origine et d'unité

Définition 2.4 On appelle changement d'origine l'opération consistant à ajouter (ou soustraire) à une quantité $a \in \mathbb{R}$ à toutes les observations

$$y_i = a + x_i, i = 1, \dots, n$$

Définition 2.5 On appelle changement d'unité l'opération consistant à multiplier (ou diviser) par une quantité $b \in \mathbb{R}$ à toutes les observations

$$y_i = bx_i, i = 1, \dots, n.$$

Définition 2.6 On appelle changement d'origine et d'unité l'opération consistant à multiplier toutes les observations par la même quantité $b \in \mathbb{R}$ puis à ajouter à la quantité $a \in \mathbb{R}$ à toutes les observations:

$$y_i = a + bx_i, i = 1, \dots, n.$$

Théorème 2.2 Si on effectue un changement d'origine et d'unité sur une variable x , alors sa moyenne est affectée du même changement d'origine et d'unité. **Démonstration** Si $y_i = a + bx_i$, alors $\bar{y} = \frac{1}{n} \sum_{i=1}^n (a + bx_i)$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = a + b \bar{x}.$$

*

Théorème 2.3 Si on effectue un changement d'origine et d'unité sur une variable x , alors sa variance est affectée par le carré du changement d'unité et pas par le changement d'origine.

Démonstration Si $y_i = a + bx_i$, alors $s_{2y}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n (a + bx_i - a - b\bar{x})^2 = b^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = b^2 s_x^2.$

$$s_x^2.$$

*

Remarque 2.71. Les paramètres de position sont tous affectés par un changement d'origine et d'unité.

3. Les paramètres de forme et d'aplatissement ne sont affectés ni par un changement d'unité ni par un changement d'origine.3. Les paramètres de forme et d'aplatissement ne sont affectés ni par un changement d'unité ni par un changement d'origine. On définit les moyennes des deux groupes :

- la moyenne du premier groupe $\bar{x} = \frac{1}{n_A} \sum_{i=1}^{n_A} x_i$,
- la moyenne du deuxième groupe $\bar{x} = \frac{1}{n_B} \sum_{i=n_A+1}^n x_i$.

La moyenne générale est une moyenne pondérée par la taille des groupes des moyennes des deux groupes. En effet $\bar{x} =$

$$\frac{1}{n} \sum_{i=1}^{n_A} x_i + \frac{1}{n} \sum_{i=n_A+1}^n x_i = \frac{1}{n} (n_A \bar{x} + n_B \bar{x}).$$

On peut également définir les variances des deux groupes:

- la variance du premier groupe $s_A^2 = \frac{1}{n_A} \sum_{i=1}^{n_A} (x_i - \bar{x})^2$,
- la variance du deuxième groupe $s_B^2 = \frac{1}{n_B} \sum_{i=n_A+1}^n (x_i - \bar{x})^2$.

Théorème 2.4 (de Huygens) La variance totale est définie par

$$s_{2x} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

se décompose de la manière suivante : $s_{2x} = n_A s_A^2 + n_B s_B^2 + n_A (\bar{x}_A - \bar{x})^2 + n_B (\bar{x}_B - \bar{x})^2$

$$\frac{1}{n} \sum_{i=1}^{n_A} (x_i - \bar{x})^2 + \frac{1}{n} \sum_{i=n_A+1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n_A} (x_i - \bar{x}_A + \bar{x}_A - \bar{x})^2 + \frac{1}{n} \sum_{i=n_A+1}^n (x_i - \bar{x}_B + \bar{x}_B - \bar{x})^2$$

Démonstration $s_{2x} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n_A} (x_i - \bar{x})^2 + \frac{1}{n} \sum_{i=n_A+1}^n (x_i - \bar{x})^2$ (2.2)

On note que $\frac{1}{n_A} \sum_{i=1}^{n_A} (x_i - \bar{x})^2 = \frac{1}{n_A} \sum_{i=1}^{n_A} (x_i - \bar{x}_A + \bar{x}_A - \bar{x})^2 = \frac{1}{n_A} \sum_{i=1}^{n_A} (x_i - \bar{x}_A)^2 + \frac{1}{n_A} \sum_{i=1}^{n_A} (\bar{x}_A - \bar{x})^2 + 2 \frac{1}{n_A} \sum_{i=1}^{n_A} (x_i - \bar{x}_A)(\bar{x}_A - \bar{x})$

2.8 Diagramme en tige et feuilles

Le diagramme en tige et feuilles ou *Stem and leaf diagram* est une manière rapide de présenter une variable quantitative. Par exemple, si l'on a la statistique ordonnée suivante:

15, 15, 16, 17, 18, 20, 21, 22, 23, 23, 23, 24, 25, 25, 26,
26, 27, 28, 28, 29, 30, 30, 32, 34, 35, 36, 39, 40, 43, 44,

la tige du diagramme sera les dizaines et les feuilles seront les unités. On obtient le graphique suivant.

The decimal point is 1 digit(s) to the right of the |

1 | 556782 | 0123334556678893 | 00245694 | 034 Ce diagramme permet d'avoir une vue synthétique de la

Évidemment, les tiges peuvent être définies par les centaines, ou des milliers, selon l'ordre de grandeur de la variable.

En langage R ## Diagramme en tige et feuilles #X=c(15,15,16,17,18,20,21,22,23,23,23,24,25,25,26,26,27

qui permet de représenter la distribution d'une variable. Ce diagramme est composé de :- Un rectangle

portant la médiane.- Ce rectangle est complété par deux segments de droites.- Pour les dessiner, on utilise les "extrêmes". **Exemple 2.9** On utilise une base de données de communes suisses de 2003 fournie par l'Office fédéral de la distribution, au sein duquel il y a beaucoup de petites communes et peu de grandes communes. Le graphique

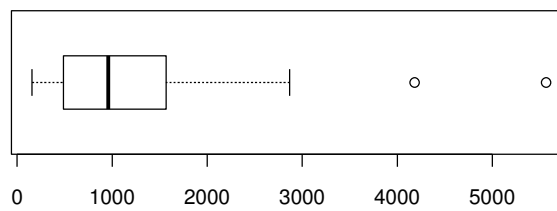


Fig. 2.5 – Boîtes à moustaches pour la variable superficie en hectares (HApoly) des communes du canton de Neuchâtel

```
# Etape 1: installation du package sampling
# dans lequel se trouve la base de données des communes belges
# choisir "sampling" dans la liste
utils::menuInstallPkgs()
# Etape 2: charge le package sampling
# choisir "sampling" dans la liste local({pkg <-
+ if(nchar(pkg)) library(pkg, character.only=TRUE)})
# Utilisation des données
data(swissmunicipalities)
attach(swissmunicipalities)
# boxplot de la sélection de
```

Exemple 2.10 On utilise une base de données belges fournie par l'Institut National (belge) de Statistique contenant des informations sur la population et les revenus des personnes physiques dans les communes. On s'intéresse à la variable "revenu moyen en euros par habitant en 2004" pour chaque commune (variable `averageincome`) et l'on aimerait comparer les 9 provinces belges : Anvers, Brabant, Flandre occidentale, Flandre orientale, Hainaut, Liège, Limbourg, Luxembourg, Namur. La Figure 2.6 contient les boîtes à moustaches des noms des provinces `b=list("Anv."=averageincome[Province==1], "Brab."=averageincome[Province==2], "Fl.`

1. Montrez que

$$s_x^2 = \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2.$$

2. Montrez que

$$s_x \leq E_t \sqrt{\frac{n-1}{2n}}.$$

3. Montrez que, si $x_i > 0$,

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| \leq 2\bar{x}.$$

Chapitre 3

Statistique descriptive bivariée

3.1 Séries statistiques bivariées

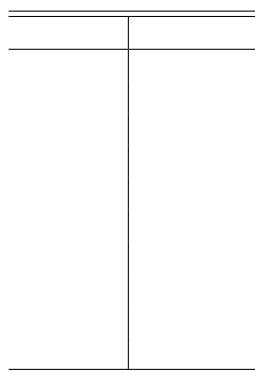
On s'intéresse à deux variables x et y . Ces deux variables sont mesurées sur n unités d'observation. Pour chaque unité, on obtient donc deux mesures. La série statistique est alors une suite de n couples des valeurs prises par les deux variables sur chaque individu : $(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)$.

Chacune des deux variables peut être, soit quantitative, soit qualitative. On examine deux cas.

- Les deux variables sont quantitatives. - Les deux variables sont qualitatives.

3.2 Deux variables

Les couples (x_i, y_i) (réels) peut toujours être représenté comme un point dans un plan $(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)$. **Exemple 3.**



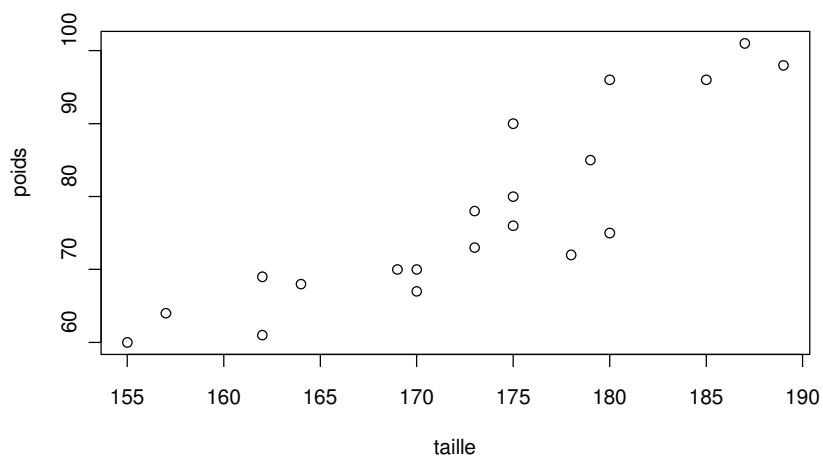


Fig. 3.1- Lenuagedepoints

3.2.2 Analyse des variables Les variables x et y peuvent être analysées séparément. On peut calculer les moyennes et les variances :

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $s_{2x} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$,
 etres dont les

$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, $s_{2y} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$. Ces paramètres sont appelés *paramètres marginaux : variances marginales*

ecarts-types marginaux, quantiles marginaux, etc... **3.2.3 Covariance** La covariance est définie $s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

Démonstration

$$\begin{aligned}
 s_{xy} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\
 &= \frac{1}{n} \sum_{i=1}^n (x_i y_i - y_i \bar{x} - \bar{y} x_i + \bar{x} \bar{y}) \\
 &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n y_i \bar{x} - \frac{1}{n} \sum_{i=1}^n \bar{y} x_i + \frac{1}{n} \sum_{i=1}^n \bar{x} \bar{y} \\
 &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} - \bar{y} \bar{x} + \bar{x} \bar{y} \\
 &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}
 \end{aligned}$$

*

3.2.4 Corrélation Le coefficient de corrélation est la covariance divisée par les deux

écarts-types marginaux :

$r_{xy} = \frac{s_{xy}}{s_x s_y}$. Le coefficient de détermination est le carré du coefficient de corrélation :

$r^2_{xy} = \frac{s^2_{xy}}{s^2_x s^2_y}$. **Remarque 3.2** – Le coefficient de corrélation mesure la dépendance linéaire entre deux

variables, mais peut cependant avoir une dépendance non-linéaire avec un coefficient de corrélation nul. **3.2.5 Droite de régression** On peut prédire y_i à partir de x_i en utilisant la droite de régression pour prédire y_i à partir de x_i . Les résidus peuvent

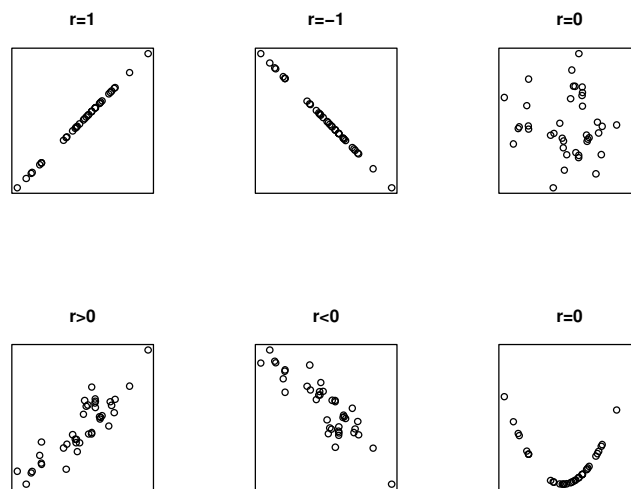


Fig. 3.2 - Exemples de nuages de points et coefficients de corrélation

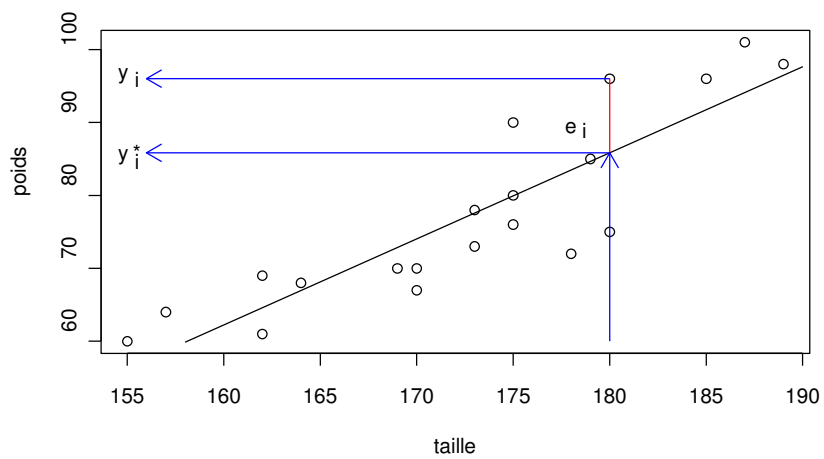


Fig. 3.3 - Le nuage de points, le résidu

En langage R

Graphique avec le r´esidus

plot(taille,poids)

segments(158,a+b*158,190,a+b*190)

segments(180,a+b*180,180,96,col="red")

#text(178,90,expression(e))

text(178.7,89.5,"i")#arrows(180,a+b*180,156,a+b*180,col="blue",length=0.14)

arrows(180,60,180,a+b*180,col="blue",length=0.14)

arrows(180,96,156,96,col="blue",length=0.14)

#text(154.8,86,expression(y))text(155.5,85.5,"i")#text(154.8,97,expression(y))text(155.5,97.8,"*")text(155.5,97.8,"*")

es qui consiste `a
chercher la droite qui minimise la somme des carr´es des r´esidus : $M(a, b) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$. **Th**

es sont donn´es par:

$b = \frac{s_{xy}}{s_{2x}}$ et $a = \bar{y} - b\bar{x}$. **D´emonstration** Le minimum $M(a, b)$ en a, b s'obtient en annulant les d´eriv´ees

et b . $\frac{\partial M(a, b)}{\partial a} = - \sum_{i=1}^n 2 (y_i - a - bx_i) = 0$ $\frac{\partial M(a, b)}{\partial b} = - \sum_{i=1}^n 2 (y_i - a - bx_i) x_i = 0$ On obtient un sys

—

—

—

—

—

—

—

—

—

—

La première équation montre que la droite passe par le point \bar{y} . On obtient

$$a = \bar{y} - b \bar{x}.$$

En remplaçant a par $\bar{y} - b \bar{x}$ dans la seconde équation, on a

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n x_i y_i - (\bar{y} - b \bar{x}) \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i^2 - b \bar{x}^2 \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{y} \bar{x} + b \bar{x}^2 - \frac{1}{n} \sum_{i=1}^n x_i^2 + \bar{x}^2 \\ &= s_{xy} - b s_x^2 \\ &= 0, \end{aligned}$$

ce qui donne

$$s_{xy} - b s_x^2 = 0.$$

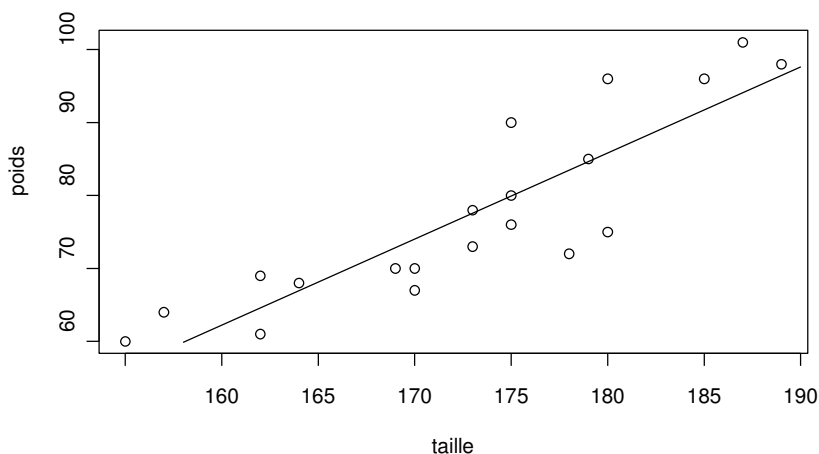
Donc $b =$

$$\frac{s_{xy}}{s_x^2}.$$

On a donc identifié les deux paramètres $b = s_{xy}/s_x^2$ (la pente) $a = \bar{y} - b \bar{x} = \bar{y} - s_{xy}/s_x^2 \bar{x}$ (la constante)

On devrait en outre vérifier qu'il s'agit bien d'un minimum en montrant que les secondes sont positives.*

La droite de régression est donc $y = a + bx = \bar{y} - s_{xy}/s_x^2 \bar{x} + s_{xy}/s_x^2 x$, ce qui peut s'écrire aussi $y - \bar{y} = s_{xy}/s_x^2 (x - \bar{x})$



Remarque 3.3 La droite de régression de y en x n'est pas la même que la droite de régression de x en y . 40

3.2.6 R esidus et valeurs ajust es

Les *valeurs ajust es* sont obtenues au moyen de la droite de r egression:

$$y_i^* = a + bx_i.$$

Les valeurs ajust es sont les 'pr edictions' des r ealis es au moyen de la variable x et de la droite de r egression de y en x . **Remarque 3.4** La moyenne des valeurs ajust es

est  gale   la moyenne des valeurs observ es. En effet,

$$\frac{1}{n} \sum_{i=1}^n y_i^* = \frac{1}{n} \sum_{i=1}^n (a + bx_i) = a + b \frac{1}{n} \sum_{i=1}^n x_i = a + b \bar{x}.$$

Or, $\bar{y} = a + b\bar{x}$, car le point (\bar{x}, \bar{y}) appartient   la droite de r egression.

Les r esidus sont les diff erences entre les valeurs observ es et les valeurs ajust es de la variable d ependante.

$$e_i = y_i - y_i^*.$$

Les r esidus repr esentent la partie inexpliqu e des y_i par la droite de r egression.

Remarque 3.5 La moyenne des r esidus est nulle. En effet $\frac{1}{n} \sum_{i=1}^n e_i = \frac{1}{n} \sum_{i=1}^n (y_i - y_i^*) = \bar{y} - \bar{y} = 0$.

- De plus, $\sum_{i=1}^n x_i e_i = 0$. La d emonstration est un peu plus difficile. **3.2.7 Sommes de carr es et var**

Définition 3.4 On appelle somme des carrés résidus (ou carré résiduel) la quantité

$$SC_{RES} = \sum_{i=1}^n e_i^2.$$

Définition 3.5 La variance résiduelle est la variance des résidus.

$$s_e^2 = \frac{SC_{RES}}{n} = \frac{1}{n} \sum_{i=1}^n e_i^2.$$

Note : Il n'est pas nécessaire de centrer les résidus sur leurs moyennes pour calculer la variance, la moyenne des résidus est nulle. **Théorème 3.3** SC_{TOT}

$$= SC_{REGR} + SC_{RES}.$$

Démonstration $SC_{TOT} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - y_{*i} + y_{*i} - \bar{y})^2 = \sum_{i=1}^n (y_i - y_{*i})^2 + \sum_{i=1}^n (y_{*i} - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - y_{*i})(y_{*i} - \bar{y})$

$$= SC_{RES} + SC_{REGR} + 2 \sum_{i=1}^n (y_i - y_{*i})(y_{*i} - \bar{y}).$$

Le troisième terme est nul. En effet, $\sum_{i=1}^n (y_i - y_{*i})(y_{*i} - \bar{y}) = \sum_{i=1}^n (y_i - y_{*i}) \sum_{i=1}^n (y_{*i} - \bar{y}) = \sum_{i=1}^n (y_i - y_{*i}) \cdot 0 = 0$.

Démonstration

$$\begin{aligned}
 s_{y^*}^2 &= \frac{1}{n} \sum_{i=1}^n (y_i^* - \bar{y})^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \left(\bar{y} + \frac{s_{xy}}{s_x^2} (x_i - \bar{x}) \right)^2 \\
 &= \frac{s_{xy}^2}{s_x^4} \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\
 &= \frac{s_{xy}^2}{s_x^2} \\
 &= s_y^2 \frac{s_{xy}^2}{s_x^2 s_y^2} \\
 &= s_y^2 r^2.
 \end{aligned}$$

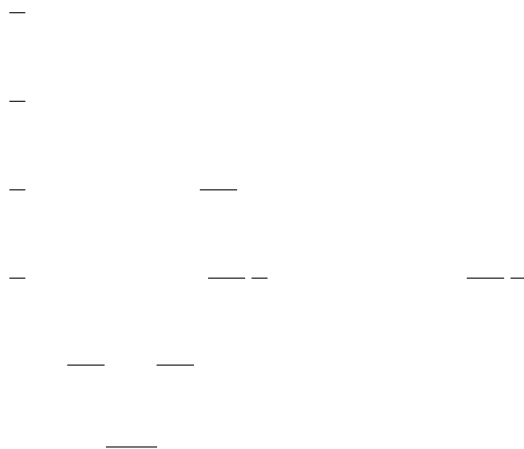
*

La variance résiduelle est la variance des résidus. $s_{e}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$

$$= \frac{1}{n} \sum_{i=1}^n e_i^2$$

Théorème 3.5 La variance résiduelle peut également s'écrire $s_e^2 = s_y^2(1 - r^2)$, où r^2 est le coefficient de détermination.

variables qualitatives 3.3.1 Données observées Si les deux variables x et y sont qualitatives, on peut les coder par des entiers.



chacune des deux variables prend comme valeurs des modalités.
Les valeurs distinctes de x et y sont notées respectivement

$$x_1, \dots, x_J, \dots, x_J$$

et

$$y_1, \dots, y_K, \dots, y_K.$$

3.3.2 Tableau de contingence

Les données observées peuvent être regroupées sous la forme d'un *tableau de contingence*

	y_1	\dots	y_k	\dots	y_K	total
x_1	n_{11}	\dots	n_{1k}	\dots	n_{1K}	$n_{1.}$
\dots	\dots		\vdots		\vdots	
x_j	n_{j1}	\dots	n_{jk}	\dots	n_{jK}	$n_{j.}$
\dots	\dots		\vdots		\vdots	
x_J	n_{J1}	\dots	n_{Jk}	\dots	n_{JK}	$n_{J.}$
total	$n_{.1}$	\dots	$n_{.k}$	\dots	$n_{.K}$	n

Les $n_{j.}$ et $n_{.k}$ sont appelés les effectifs marginaux. Dans ce tableau,

- $n_{j.}$ représente le nombre de fois que la modalité x_j apparaît,
- $n_{.k}$ représente le nombre de fois que la modalité y_k apparaît, - n_{jk} représente le nombre de fois que

On a les relations $\sum_{j=1}^J n_{jk} = n_{.k}$, pour tout $k = 1, \dots, K$, $\sum_{k=1}^K n_{jk} = n_{j.}$, pour tout $j = 1, \dots, J$, et $\sum_{j=1}^J \sum_{k=1}^K n_{jk} = n$.

Le Tableau 3.1 reprend le tableau de contingence. Tab. 3.1 – Tableau des effectifs n_{jk}

—

—

Le tableaude fréquences est

	y_1	\cdots	y_k	\cdots	y_K	total
x_1	f_{11}	\cdots	f_{1k}	\cdots	f_{1K}	$f_{1.}$
\vdots	\vdots		\vdots		\vdots	
x_j	f_{j1}	\cdots	f_{jk}	\cdots	f_{jK}	$f_{j.}$
\vdots	\vdots		\vdots		\vdots	
x_J	f_{J1}	\cdots	f_{Jk}	\cdots	f_{JK}	$f_{J.}$
total	$f_{.1}$	\cdots	$f_{.k}$		$f_{.K}$	1

Exemple 3.3 Le Tableau 3.2 reprend le tableau des séquences.

Tab. 3.2- Tableau des fréquences

	Bleu	Vert	Marron	Total
Homme	0.05	0.25	0.10	0.40
Femme	0.10	0.30	0.20	0.60
Total	0.15	0.55	0.30	1.00

3.3.4 Profils lignes et profils colonnes

en colonnes (appelées aussi *profils lignes* et *profils colonnes*). Les profils lignes sont définis par $f_{(j)k} = n_{jk}n_{j.} =$

Tab. 3.4 – Tableaux des profils colonnes

	Bleu	Vert	Marron	Total
Homme	0.33	0.45	0.33	0.40
Femme	0.67	0.55	0.67	0.60
Total	1.00	1.00	1.00	1.00

3.3.5 Effectifs théoriques et khi-carré

On cherche souvent une interaction entre des lignes et des colonnes. Pour mettre en évidence ce lien, on construit un tableau d'effectifs qui représente la situation où les variables ne sont pas liées (indépendance). Ces effectifs sont construits de la manière suivante :

$$n_{jk}^* = \frac{n_{j.} n_{.k}}{n}.$$

Les effectifs observés n_{jk} ont les mêmes marges que les effectifs théoriques. Enfin, les écarts à l'indépendance sont définis par $e_{jk} = n_{jk} - n_{jk}^*$.

j, k .

- La dépendance du tableau se mesure au moyen du khi-carré

$$\chi^2_{obs} = \sum_{k=1}^K \sum_{j=1}^J \frac{(n_{jk} - n_{jk}^*)^2}{n_{jk}^*} = \sum_{j=1}^J \sum_{k=1}^K \frac{e_{jk}^2}{n_{jk}^*}. \quad (3.1)$$

- Le khi-carré peut être normalisé pour ne plus dépendre du nombre d'observations. On définit le phi-deux par : $\phi^2 = \chi^2_{obs} / n$. Le ϕ^2 ne dépend plus du nombre d'observations. Il est possible de montrer que $\phi^2 \leq 1$ (tableau ait le même nombre de lignes que de colonnes). **Exemple 3.5** Le Tableau 3.5 reprend le tableau de

Tab. 3.6 – Tableau des écarts à l'indépendance

	Bleu	Vert	Marron	Total
Homme	-2	6	-4	0
Femme	2	-6	4	0
Total	0	0	0	0

Tab. 3.7 – Tableaudes e_{jk}^2/n_{jk}^*

	Bleu	Vert	Marron	Total
Homme	0.33	0.82	0.67	1.82
Femme	0.22	0.55	0.44	1.21
Total	0.56	1.36	1.11	3.03

- Le khi-carré observé vaut $\chi_{2obs}^2 = 3.03$. - Le phi-deux vaut $\phi_2 = 0.01515$. - Comme le tableau a deux lignes et deux colonnes, la valeur de ϕ_2 est égale à ϕ . - On a $V = 0.01515$. La dépendance entre les deux variables est très faible.

En langage R yeux = c(rep("bleu", times = 10), rep("vert", times = 50), rep("marron", times = 20),

rep("bleu", times = 20), rep("vert", times = 60), rep("marron", times = 40))
 sexe = c(rep("homme", times = 80), rep("femme", times = 80))
 Tableau de contingence: effectifs n_{jk} Niveau d'instruction Statut professionnel du fils du fils par rapport

Tab. 3.9 – Tableau des fréquences f_{jk}

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	0.322	0.231	0.147	0.700
égal	0.055	0.079	0.058	0.192
inférieur	0.017	0.038	0.053	0.108
total	0.394	0.349	0.257	1.000

Tab. 3.10 – Tableau des profils lignes

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	0.460	0.330	0.210	1
égal	0.288	0.413	0.300	1
inférieur	0.156	0.356	0.489	1
total	0.394	0.349	0.257	1

Tab. 3.11 – Tableau des profils colonnes

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	0.817	0.662	0.570	0.700
égal	0.140	0.228	0.224	0.192
inférieur	0.043	0.110	0.206	0.108
total	1.111			

Tab. 3.12 – Tableau des effectifs théoriques n_{*jk}

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	114.72	114.72	114.72	344.16
égal	12.52	12.52	12.52	37.56
inférieur	1.24	1.24	1.24	3.72
total	128	128	128	384

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	114.72	114.72	114.72	344.16
égal	12.52	12.52	12.52	37.56
inférieur	1.24	1.24	1.24	3.72
total	128	128	128	384

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	114.72	114.72	114.72	344.16
égal	12.52	12.52	12.52	37.56
inférieur	1.24	1.24	1.24	3.72
total	128	128	128	384

X\Y	Plus élevée	Egal	inférieur	total
plus élevée	114.72	114.72	114.72	344.16
égal	12.52	12.52	12.52	37.56
inférieur	1.24	1.24	1.24	3.72
total	128	128	128	384

Exercices

Exercice 3.1 La consommation de crèmes glacées

Les crèmes glacées par individu ont été mesurées pendant 30 jours. L'objectif est de déterminer si la consommation de crèmes glacées dépend de la température. Les données sont dans le tableau 3.15. On sait en outre que Tab. 3.15–Consommation de crèmes glacées

Tab. 3.15–Consommation de crèmes glacées

consommation	température	consommation	température	consommation	température
386	41	286	28	319	44
374	56	298	26	307	40
393	63	329	32	284	32
425	68	318	40	326	27
406	69	381	55	309	28
344	65	381	63	359	33
327	61	470	72	376	41
288	47	443	72		52
269	32	386	67		64
256	24	342	60		71

$\sum_{i=1}^n y_i = 10783$, $\sum_{i=1}^n x_i = 1473$, $\sum_{i=1}^n y_i^2 = 4001293$, $\sum_{i=1}^n x_i^2 = 80145$, $\sum_{i=1}^n x_i y_i = 553747$, 1. Donnez les moyennes

2. Donnez la droite de régression, avec comme variable dépendante la consommation de glaces et comme variable explicative la température. 3. Donnez la valeur ajustée et le résidu pour la première observation du tableau. 4. Déterminez le coefficient de détermination R^2 . 5. Établissez, sur la base du modèle, une

Exercice 3.3 Considérons un échantillon de 10 fonctionnaires (ayant entre 40 et 50 ans) d'un minist`re. Soit X le nombre d'années de service et Y le nombre de jours d'absence pour raison de maladie (au cours de l'année précédente) déterminé pour chaque personne appartenant à l'échantillon.

x_i	2	14	16	8	13	20	24	7	5	11
y_i	3	13	17	12	10	8	20	7	2	8

1. Représentez le nuage de points.
2. Calculez le coefficient de corrélation entre X et Y .
3. Déterminez l'équation de la droite de régression de Y en fonction de X .
4. Déterminez la qualité de cet ajustement.
5. Établissez, sur base de ce modèle, le nombre de jours d'absence pour un fonctionnaire ayant 22 ans de service. 50

Chapitre 4

Théorie des indices, mesures d'inégalité

4.1 Nombres indices

4.2 Définition

Un indice est la

référence. Prenons l'exemple du tableau 4.1 contenant le prix (fictif) d'un bien de consommation de 2000 à 2006. Le temps varie de 0, 1, 2, ..., 6 et est considéré comme le temps de référence par rapport auquel l'indice est calculé.

Tab. 4.1- Tableau du prix d'un bien de consommation de 2000 à 2006

année t	prix p_t
2000 0	2.00
2001 1	2.30
2002 2	2.40
2003 3	2.80
2004 4	3.00
2005 5	3.50
2006 6	4.00

L'indice

142.86	66.67	76.67	80.00	93.33	100.00	116.67	133.33	157.14	165.71	168.57	80.00	85.71	100.00	114.29	5
--------	-------	-------	-------	-------	--------	--------	--------	--------	--------	--------	-------	-------	--------	--------	---

4.2.1 Propriétés des indices

Considérons un indice quelconque $I(t/0)$. On dit que cet indice possède les propriétés de

- *réversibilité* si $I(t/0) = 100 \times \frac{1}{I(0/t)}$,
- *identité* si $I(t/t) = 100$,
- *circularité* (ou *transitivité*) si $I(t/u) \times I(u/v) = 100 \times I(t/v)$.

Il est facile de montrer que ces quatre propriétés sont satisfaites pour un indice simple.

4.2.2 Indices synthétiques

Quand on veut calculer un indice à partir de plusieurs prix, problème devient sensiblement plus compliqué. Un indice synthétique est une grandeur d'un ensemble de biens par rapport à une année de référence. On ne peut pas construire un indice synthétique en additionnant simplement des indices simples. Il faut, en effet, tenir compte des quantités achetées.

Pour calculer un indice de prix de n biens de consommation numérotés de 1, 2, .., n , on utilise la notation suivante: - p_{ti} représente le prix du bien de consommation i au temps t ,

- q_{ti} représente la quantité de biens i consommée au temps t .

Considérons par exemple le Tableau 4.3 qui contient 3 biens de consommation et pour lesquels on connaît les prix et les quantités achetées. Tab. 4.3 – Exemple : prix et quantités de trois biens pendant 3 ans

Temps	0	1	2			
	Prix (p_{0i})	Quantités (q_{0i})	Prix (p_{1i})	Quantités (q_{1i})	Prix (p_{2i})	Quantités (q_{2i})
Bien 1	100	14	150	10	200	8
Bien 2	60	10	50	12	40	14
Bien 3	160	4	140	5	140	5

Il existe deux méthodes fondamentales pour calculer les indices de prix, l'indice de Paasche et l'indice de Laspeyres. **4.2.3 Indice de Laspeyres** L'indice de Laspeyres, est défini par $L(t/0) = 100 \times \frac{\sum_{i=1}^n q_{0i} p_{ti}}{\sum_{i=1}^n q_{0i} p_{0i}}$. L'indice de Laspeyres ne possède

Exemple 4.1 Si on utilise les données du tableau 4.3, les indices de Laspeyres sont les suivants

$$L(1/0) = 100 \times \frac{\sum_{i=1}^n q_{0i} p_{1i}}{\sum_{i=1}^n q_{0i} p_{0i}} = 100 \times \frac{14 \times 150 + 10 \times 50 + 4 \times 140}{14 \times 100 + 10 \times 60 + 4 \times 160} = 119.6970,$$

$$L(2/0) = 100 \times \frac{\sum_{i=1}^n q_{0i} p_{2i}}{\sum_{i=1}^n q_{0i} p_{0i}} = 100 \times \frac{14 \times 200 + 10 \times 40 + 4 \times 140}{14 \times 100 + 10 \times 60 + 4 \times 160} = 142.4242,$$

$$L(2/1) = 100 \times \frac{\sum_{i=1}^n q_{1i} p_{2i}}{\sum_{i=1}^n q_{1i} p_{1i}} = 100 \times \frac{10 \times 200 + 12 \times 40 + 5 \times 140}{10 \times 150 + 12 \times 50 + 5 \times 140} = 113.5714.$$

4.2.4 Indice de Paasche L'indice de Paasche, est défini par $P(t/0) = 100 \times$

$$\frac{\sum_{i=1}^n q_{ti} p_{ti}}{\sum_{i=1}^n q_{ti} p_{0i}}.$$

On utilise, pour le calculer, les quantités q_{ti} du temps par rapport auquel on veut calculer l'indice.

L'indice de Paasche peut aussi être présenté comme une moyenne harmonique des indices simples. Soient l'indice simple du bien i : $I_i(t/0) = 100 \times \frac{p_{ti}}{p_{0i}}$, et le poids w_{ti} correspondant à la recette totale

—

des recettes au temps t : $P(t/0) = \frac{\sum_{i=1}^n w_{ti}}{\sum_{i=1}^n w_{ti} / I_i(t/0)} = \frac{\sum_{i=1}^n p_{ti} q_{ti}}{\sum_{i=1}^n p_{ti} q_{ti} / (100 \times \frac{p_{ti}}{p_{0i}})} = 100 \times \frac{\sum_{i=1}^n q_{ti} p_{0i}}{\sum_{i=1}^n q_{ti} p_{ti}}$.

— — —

difficile à calculer que l'indice de Laspeyres, car on doit connaître les quantités pour chaque valeur de t .

Fisher L'indice de Laspeyres est en général plus grand que l'indice de Paasche, ce qui peut s'expliquer par

— — —
— — —
— — —

Fisher a proposé d'utiliser un compromis entre l'indice de Paasche et de Laspeyres en calculant simplement la moyenne géométrique de ces deux indices

$$F(t/0) = \sqrt{L(t/0) \times P(t/0)}.$$

L'avantage de l'indice de Fisher est qu'il jouit de la propriété de réversibilité.

Exemple 4.3 Si on utilise toujours les données du tableau 4.3, les indices de Fisher sont les suivants:

$$F(1/0) = \sqrt{L(1/0) \times P(1/0)} = 115.3242,$$

$$F(2/0) = \sqrt{L(2/0) \times P(2/0)} = 129.2052,$$

$$F(2/1) = \sqrt{L(2/1) \times P(2/1)} = 111.7715.$$

4.2.6 L'indice de Sidgwick L'indice de Sidgwick est la moyenne arithmétique des indices de Paasche

$$S(t/0) = \frac{L(t/0) + P(t/0)}{2}.$$

4.2.7 Indices chaînés Le défaut principal des indices de Laspeyres, de Paasche, de Fisher et de Sidgwick

est qu'ils ne possèdent pas la propriété de circularité. Un indice qui possède cette propriété est appelé un indice chaîné. Pour construire un indice chaîné, avec l'indice de Laspeyres, on peut faire un produit d'indices de Laspeyres annuels.

$$CL(t/0) = 100 \times \frac{L(t/t-1)}{100} \times \frac{L(t-1/t-2)}{100} \times \dots \times \frac{L(2/1)}{100} \times \frac{L(1/0)}{100}.$$

Pour calculer un tel indice, on doit évidemment connaître les quantités pour chaque année de t . L'indice suisse des prix à la consommation est un indice chaîné de Laspeyres. **Exemple 4.4** En utilisant encore les

suivants : $CL(1/0) = L(1/0) = 119.6970$, $CL(2/1) = L(2/1) = 113.5714$, $CL(2/0) = L(2/1) \times L(1/0) / 100 = 135.9416$.

4.3.2 Courbe de Lorenz

Plusieurs indices d'égalié sont liés à la courbe de Lorenz. On note

$$x_1, \dots, x_{(1)}, \dots, x_{(n)}$$

les revenus des n individus de la population. On note également

$$x_{(1)}, \dots, x_{(i)}, \dots, x_{(n)},$$

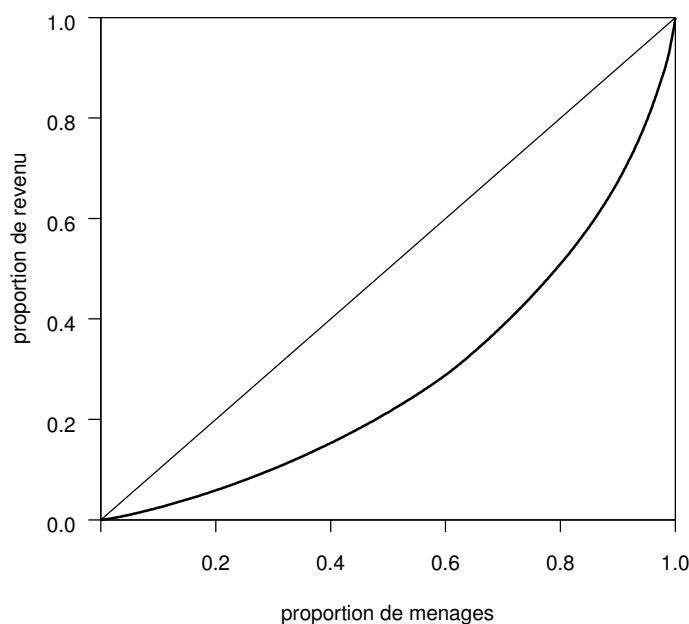
la statistique d'ordre, c'est-à-dire les revenus triés par ordre croissant.

Notons maintenant q_i la proportion de revenus par rapport au revenu total que les individus ayant les plus bas revenus, ce point

$$q_i = \frac{\sum_{j=1}^i x_{(j)}}{\sum_{j=1}^n x_{(j)}} \text{ avec } q_0 = 0 \text{ et } q_n = 1.$$

La courbe de Lorenz est la représentation graphique de la fonction qui à la part des individus les moins riches associe la part y du revenu total qu'ils perçoivent. Plus précisément la courbe de Lorenz relie les points $(i/n, q_i)$ pour $i = 1, \dots, n$. En abscisse, on a donc une proportion d'individus classés par ordre de revenu, et en ordonnée la proportion du revenu total reçu par ces individus.

Exemple 4.5 On utilise une enquête ménagerie sur le revenu des ménages des Philippines appelée cos. Cette enquête de 1997 sur le revenu des ménages a été produite par l'Office philippin de Statistique. La courbe de Lorenz est présentée en Figure 4.1. Fig. 4.1 - Courbe de Lorenz



Remarque 4.1 Sur le graphique, on indique toujours la diagonale. La courbe de Lorenz est égale à la diagonale si et seulement si la répartition des revenus est parfaitement égale.


```
## Courbe de Lorenz et indices d'in egalit e
```

```
## Etape 1: on installe le package ineq
```

```
utils::menuInstallPkgs()
```

```
# choisir 'ineq' dans la liste
```

```
## Etape 2: on charge le package ineq
```

```
local({pkg<-select.list(sort(.packages(all.available=TRUE)))
```

```
+ if(nchar(pkg))library(pkg,character.only=TRUE)})
```

```
# choisir 'ineq' dans la liste
```

```
## Utilisation de la base de donn ees llocos
```

```
# Enqu etes sur le revenu de l'Office de Statistique Philippin
```

```
data(llocos)attach(llocos)# plot(Lc(income),xlab="proportion de m enages",
```

```
ylab="proportion de revenu",main="")4.3.3 Indice de Gini L'indice de Gini, not e  $G$ , est  gal   deux fois la surface
```

Il est possible de montrer que: $G = \frac{1}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=1}^n |x_i - x_j|$. En utilisant la statistique d'ordre $x_{(1)}, \dots, x_{(i)}, \dots, x_{(n)}$, l'indice de

 crire

$G = 1 - \frac{2}{n} \sum_{i=1}^n i x_{(i)} / \sum_{i=1}^n x_{(i)}$. L'indice de Gini est compris entre 0 et 1. S'il est proche de 0, tous les revenus sont  gaux. S'il est proche de 1, les revenus sont tr es in gaux.

4.3.4 Indice de Hoover L'indice d' equit e de partition de Hoover (ou *Robin Hood*) est  gal   la courbe de Lorenz, car il est possible de montrer qu'il correspond   la plus grande distance verticale entre la courbe de Lorenz et la diagonale.

Le quintile shareratio est défini par

$$QSR = \frac{S_{80}}{S_{20}}.$$

Le decile shareratio est défini par

$$DSR = \frac{S_{90}}{S_{10}}.$$

Ces quantités sont toujours plus grandes que 1 et augmentent avec l'inégalité. Ces deux rapports sont facilement interprétables, par exemple, si $QSR = 5$, cela signifie que le revenu moyen des 20% les plus riches est 5 fois plus grand que le revenu moyen des 20% les plus pauvres.

4.3.6 Indice de pauvreté Un indice simple de pauvreté consiste à calculer le pourcentage de la population

moitié de la médiane. **4.3.7 Indices selon les pays** Le tableau 4.4 reprend pour tous les pays l'indice de Gini

et le rapport des 20% les plus riches sur les 20% les plus pauvres. (référence: United Nations 2005 Development Programme Report, page 270).

Exercices **Exercice 4.1** Étudiez les propriétés (circularité, réversibilité, identité

et transitivité) de tous les indices de prix présentés. 57

Tab.4.4-Mesuresdel'in´egalit  parpays

Rang	Pays	Indice deGini	DSR	QSR	Ann��e del'enqu��te
1	Denmark	24.7	8.1	4.3	1997
2	Japan	24.9	4.5	3.4	1993
3	Sweden	25	6.2	4	2000
4	Belgium	25	7.8	4.5	1996
5	CzechRepublic	25.4	5.2	3.5	1996
6	Norway	25.8	6.1	3.9	2000
7	Slovakia25.8		6.7	4	1996
8	BosniaandHerzegovina	26.2	5.4	3.8	2001
9	Uzbekistan26.8		6.1	4	2000
10	Finland26.9		5.6	3.8	2000
11	Hungary26.9		5.5	3.8	2002
12	RepublicofMacedonia	28.2	6.8	4.4	1998
13	Albania28.2		5.9	4.1	2002
14	Germany28.3		6.9	4.3	2000
15	Slovenia28.4		5.9	3.9	1998
16	Rwanda28.9		5.8	4	1983
17	Croatia29		7.3	4.8	2001
18	Ukraine296.4			4.3	1999
19	Austria307.6			4.7	1997
20	Ethiopia306.6			4.3	1999
21	Romania30.38.1			5.2	2002
22	Mongolia30.317.89.1				1998
23	Belarus30.46.94.6				2000
24	Netherlands30.99.25.1				1999
25	Russia317.14.8				2002
26	SouthKorea31.67.84.71998				
27	Bangladesh31.86.84.62000				
28	Lithuania31.97.95.12000				
29	Bulgaria31.99.95.82001				
30	Kazakhstan32.37.55.12003				
31	Spain32.595.4199032	India32.57.34.9199933	Tajikistan32.67.85.2200334	France32.79.15.6199535	Pakistan337.64.819983
a57.833.117.92000117	Brazil59.36826.42001118	Guatemala59.955.124.42000119	Swaziland60.949.723.81994120	CentralA	

Chapitre 5

Séries temporelles, filtres, moyennes mobiles et désaisonnalisation

5.1 Définitions générales et exemples

5.1.1 Définitions **Définition 5.1** Une série temporelle est une suite d'observations d'une quantité mesurée à des instants t dans le temps.

On énonce en général l'hypothèse que les intervalles de temps sont égaux, c'est-à-dire que la série est notée $y_1, \dots, y_t, \dots, y_T$. On note également $T = \{1, 2, \dots, t, \dots, T\}$ l'ensemble des instants auxquels les observations sont réalisées.

Une série temporelle peut se composer de :
- une tendance T_t ,
- une composante cyclique C_t (nous n'étudierons pas la saisonnalité),
- une composante saisonnière, et désaisonnaliser la série, c'est réaliser une prévision pour des valeurs de t à partir de données historiques.
QTR: Quarter, trimestres depuis le 1er trimestre 1978 jusqu'au 4^{ème} trimestre 1985

- DISH:Nombre de lave-vaisselles(dishwashers)expédés(milliers)
- DISP:Nombre de broyeur d'ordures(disposers)expédés(milliers)
- FRIG:Nombre de réfrigérateursexpédés(milliers)
- WASH:Nombre de machines à laver(washing machine)expédés(milliers)
- DUR:Dépenses en biens durables USA(milliards de dollars de 1982)
- RES:Investissement résidentiel privé USA(milliards de dollars de 1982)

Tab.5.1-Biens manufacturés aux USA

QTR	DISH	DISP	FRIG	WASH	DUR	RES
1	1841	798	1317	1271	252.6	172.9
2	2957	837	1615	1295	272.4	179.8
3	3999	821	1662	1313	270.9	180.8
4	4960	858	1295	1150	273.9	178.6
5	5894	837	1271	1289	268.9	174.6
6	6851	838	1555	1245	262.9	172.4
7	7863	832	1639	1270	270.9	170.6
8	8878	818	1238	1103	263.4	165.7
9	9792	868	1277	1273	260.6	154.9
10	1058	962	312	581	231.9	124.1
11	1165	766	214	171	242.7	126.8
12	1269	982	211	851	101	142.2
13	1367	587	111	961	181	1258.7
14	1465	279	114	1011	162	248.4
15	1562	875	914	1711	902	55.1
16	1652	973	491	911	252	40.4
17	1748	070	694	310	362	47.7
18	1853	058	211	751	019	249.1
19	1103	419	557	659	126	910
20	472	51	810	012	060	283
21	797	391	826	201	158	216
22	115	821	658	867	110	210

1853058211751019249.1103.41955765912691047251.8100.120602837973918262.0115.8216588671102
 'érateurs vendus a manifestement une composante saisonnière et une tendance. **En langage R** QTR=c(1,2,3,4,5,

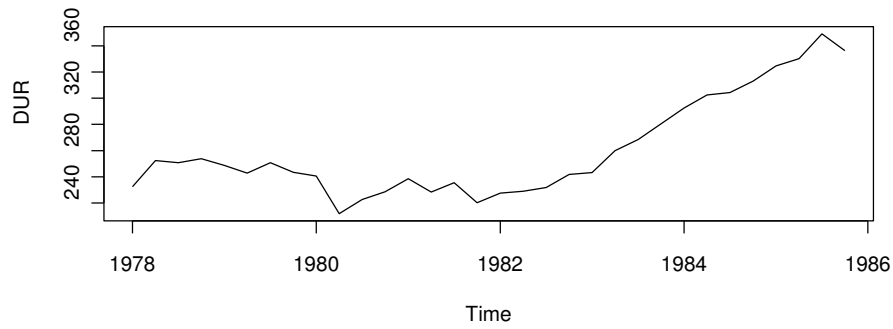


Fig.5.1-D ´epenses en biens durables USA (milliards de dollars de 1982)

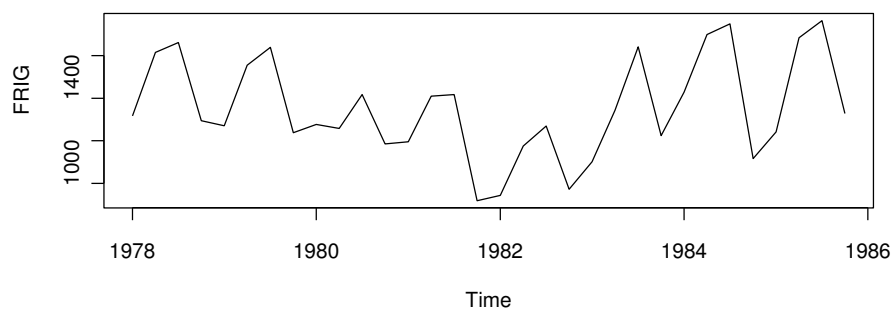


Fig.5.2-Nombre de réfrigérateurs vendus de 1978 à 1985

```

DISP=c(798,837,821,858,837,838,832,818,868,623,662,822,871,791,759,734,706,
582,659,837,867,860,918,1017,1063,955,973,1096,1086,990,1028,1003)
FRIG=c(1317,1615,1662,1295,1271,1555,1639,1238,1277,1258,1417,1185,1196,1410,1417,919,943,1175,
1242,1684,1764,1328)
WASH=c(1271,1295,1313,1150,1289,1245,1270,1103,1273,1031,1143,1101,1181,1144,161.9,159.9,170.5,173.1,170.3,169.6,170.3,172.9,175,179.4)
plot(QTR,DUR,type="l")
plot(QTR,FRIG,type="l")

```

Tab.5.2-Indicedesprix`alaconsommation(France)

p_t	1970	1971	1972	1973	1974	1975	1976	1977	1978
janvier	97.9	102.5	108.3	115.5	127.4	145.9	159.9	174.3	190.3
f`evrier98.2		103.0	108.9	115.8	129.1	147.0	161.0	175.5	191.7
mars98.5		103.4	109.4	116.4	130.6	148.2	162.4	177.1	193.4
avril99.0		104.0	109.8	117.2	132.7	149.5	163.8	179.4	195.5
mai99.4		104.7	110.4	118.3	134.3	150.6	164.9	181.1	197.4
juin99.8105.1			111.0	119.2	135.8	151.7	165.6	182.5	198.9
juillet100.0105.6			111.9	120.2	137.5	152.8	167.2	184.1	201.5
ao`ut100.4106.0			112.5	121.0	138.6	153.8	168.4	185.1	202.5
septembre100.8106.5			113.2	122.1	140.1	155.1	170.2	186.7	203.8
octobre101.2107.1114.2				123.4	141.8	156.3	171.8	188.2	205.7
novembre101.6107.5114.9				124.5	143.1	157.3	173.2	188.9	206.8
d`ecembre101.9108.0115.5125.3					144.3	158.2	173.8	189.4	207.8

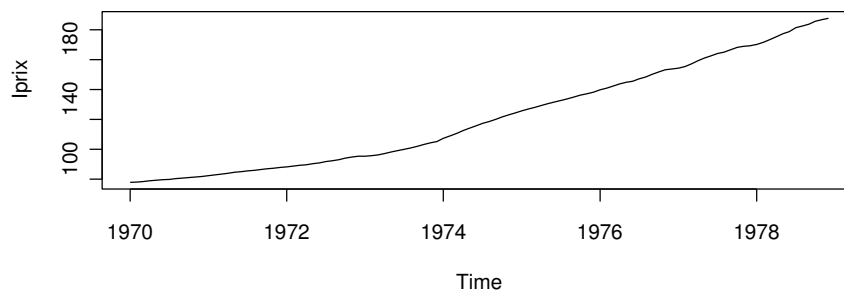


Fig.5.3-Indicedesprix`alaconsommation p_t Time

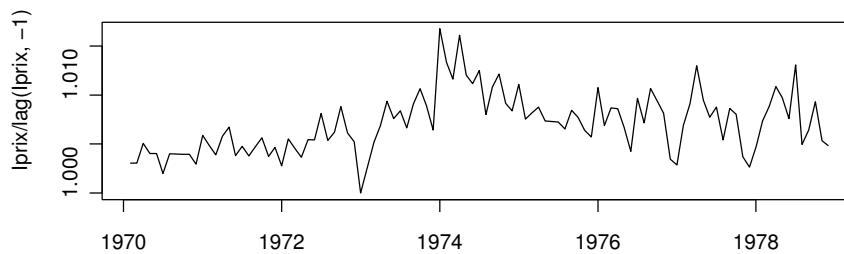


Fig.5.4-Rapportmensuellesindicesdeprix p_t/p_{t-1} Time

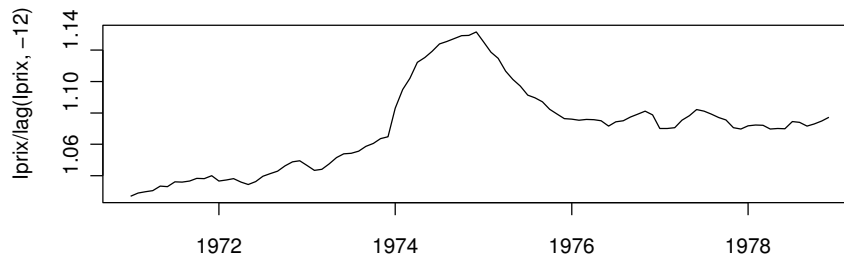


Fig.5.5-Rapportenglisementannuellesindicesdeprix p_t/p_{t-12} #Iprix=c(97.9,98.2,98.5,99,99.4,99.8,100,100.4,100.8,101.2,101.6,101.9)

108.3,108.9,109.4,109.8,110.4,111,111.9,112.5,113.2,114.2,114.9,115.5,
 115.5,115.8,116.4,117.2,118.3,119.2,120.2,121,122.1,123.4,124.5,125.3,
 127.4,129.1,130.6,132.7,134.3,135.8,137.5,138.6,140.1,141.8,143.1,144.3,
 145.9,147,148.2,149.5,150.6,151.7,152.8,153.8,155.1,156.3,157.3,158.2,
 159.9,161,162.4,163.8,164.9,165.6,167.2,168.4,170.2,171.8,173.2,173.8,
 174.3,175.5,177.1,179.4,181.1,182.5,184.1,185.1,186.7,188.2,188.9,189.4,
 190.3,191.7,193.4,195.5,197.4,198.9,201.5,202.5,203.8,205.7,206.8,207.8)

lprix<-ts(lprix,start=c(1970,1),frequency=12)

plot(lprix)plot(lprix/lag(lprix,-1))plot(lprix/lag(lprix,-12))

Exemple 5.4 Données du nombre de voyageurs-kilom

etres en millions de
 kilom`etres. Tab. 5.3–Trafic du nombre de voyageurs SNCF

mois/ann´ee	janv.	fév.	mars	avril	mai	juin	juil.	août	sept.	oct.	nov.	déc.
1963	1750	1560	1820	2090	1910	2410	3140	2850	2090	1850	1630	2420
1964	1710	1600	1800	2120	2100	2460	3200	2960	2190	1870	1770	2270
1965	1670	1640	1770	2190	2020	2610	3190		2140	1870	1760	2360
1966	1810	1640	1860	1990	2110	2500	3030		2160	1940	1750	2330
1967	1850	1590	1880	2210	2110	2480	2880		2100	1920	1670	2520
1968	1834	1792	1860	2138	2115	2485	2581			1936	1784	2391
1969	1798	1850	1981	2085	2120	2491	2834			2085	1856	2553
1970	1854	1823	2005	2418	2219	2722	2912				1910	2537
1971	2008	1835	2120	2304	2264	2175	2928				2009	2546
1972	2084	2034	2152	2252	2231	1826	8429				1978	2723
1973	2081	2112	2279	2661	2281	2929	3089					2862
1974	2232	2248	2421	2710	2505	3021	3327					2876
1975	2481	2428	2596	2923	2795	3287	3598					3266
1976	2667	2668	2804	2806	2976	3430	3705					
1977	2706	2586	2796	2978	3053	3463	3649					
1978	2820	2857	3306	3333	3314	1351	2374					

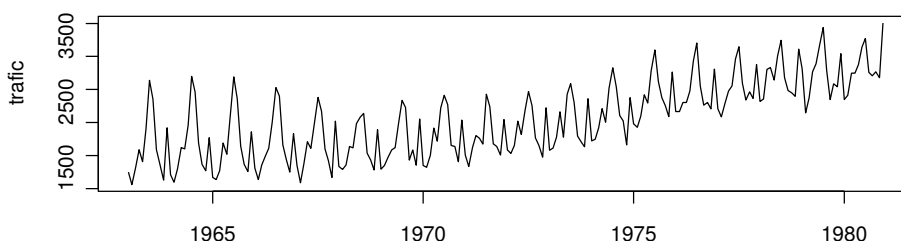


Fig. 5.6–Trafic du nombre de voyageurs SNCF

5.2 Description de la tendance

5.2.1 Les principaux modèles

Plusieurs types de modèles peuvent être utilisés pour décrire la tendance.

- Modèles dépendant du temps. La série dépend directement du temps. Le modèle peut être additif:

$$y_t = f(t) + E_t,$$

ou multiplicatif

$$y_t = f(t) \times E_t.$$

- Modèles explicatifs statiques: la série chronologique dépend des valeurs prises par une ou plusieurs autres séries chronologiques. y_t

$$y_t = f(x_t) + E_t$$

Le cas linéaire est le plus facile à traiter $y_t = b_0 + b_1 x_t + E_t$.

$$y_t = b_0 + b_1 x_t + E_t.$$

- Modèles auto-projectifs. La série chronologique au temps t dépend de ses propres valeurs passées

$$y_t = f(y_{t-1}, y_{t-2}, y_{t-3}, \dots, y_{t-p}) + E_t$$

- Modèles explicatifs dynamiques: la série chronologique dépend des valeurs présentes et passées d'une ou de plusieurs autres séries chronologiques, par exemple: $y_t = \mu + \theta_1 y_{t-1} + \theta_2 y_{t-2} + \dots + \theta_p y_{t-p} + \varphi_1 x_{t-1} + \varphi_2 x_{t-2}$

$$+ \dots + \varphi_q x_{t-q} + E_t.$$

5.2.2 Tendance linéaire

La tendance la plus simple est la linéaire. On peut estimer les paramètres au moyen

des méthodes de

moindres carrés. C'est une régression simple. $T_t = a + bt$.

5.2.3 Tendance quadratique

On peut utiliser une

des moindres carrés. C'est une régression avec deux variables explicatives. $T_t = a + bt + ct^2$.

5.2.4 Tendance exponentielle

linéaire, on peut se ramener à un problème linéaire. En posant $z_t = 1/T_t$, on a $z_t = 1 + be^{-at}$.

$$\begin{aligned}
 z_{t+1} &= \frac{1+be^{-a(t+1)}}{c} \\
 &= \frac{1+be^{-at}e^{-a}}{c} \\
 &= \frac{1+(1+be^{-at})e^{-a}-e^{-a}}{c} \\
 &= \frac{1-e^{-a}}{c} + z_t e^{-a}.
 \end{aligned}$$

En posant $\alpha =$

$$\frac{1-e^{-a}}{c}, \text{ et } \beta = e^{-a}.$$

on obtient $z_{t+1} = \alpha + \beta z_t$,

ce qui est un modèle de type auto-projectif. On peut alors déterminer les valeurs de α et β par une simple régression linéaire. Ensuite on déduit a de la manière suivante :

$$a = -\log \beta,$$

et comme $\alpha = 1 - e^{-a}c = 1 - \beta$

$$c = \frac{\alpha}{1 - \beta},$$

on détermine c par $c = \frac{1 - \beta}{\alpha}$. Enfin, on remarque que $z_t - 1/c = be^{-at}c$, on peut donc déterminer autant de valeurs de b que l'on veut.

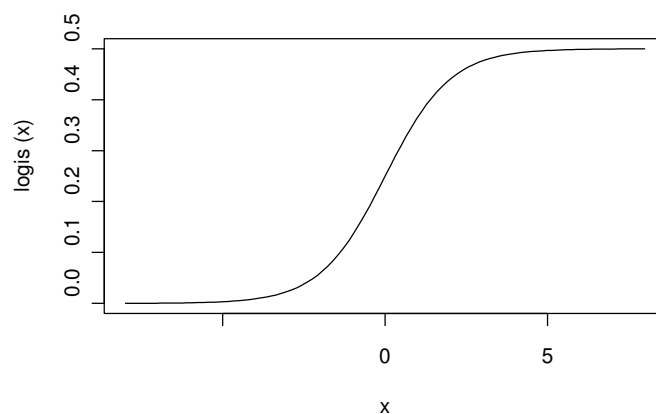


Fig.5.7-Exemple de fonction logistique avec $c=0.565$

5.3 Opérateurs de décalage et de différence

5.3.1 Opérateurs de décalage

Afin de simplifier la notation, on utilise des opérateurs de décalage. On définit l'opérateur de décalage "retard" (en anglais *lag operator*) L par

$$Ly_t = y_{t-1},$$

et l'opérateur (en anglais *forward operator*) "avance" F

$$Fy_t = y_{t+1},$$

l'opérateur identité I par

$$Iy_t = y_t.$$

L'opérateur avance est l'inverse de l'opérateur retard

$$FL = LF = I.$$

On peut donc écrire $F^{-1} = L$ et

$$L^{-1} = F.$$

On a également $L^2 y_t = LLy_t = y_{t-2}$, $-L^q y_t = y_{t-q}$, $-F^q y_t = y_{t+q}$, $-L^0 = F^0 = I$, $-L^{-q} y_t = F^q y_t = y_{t+q}$.
5.3.2 Opérateur différentiel dans le graphique 5.9. **En langage R** ## Tendances linéaire et différence $\# \text{lin} = 10 + 0.3 * (0:50) + \text{rnorm}(50, 0, 1)$

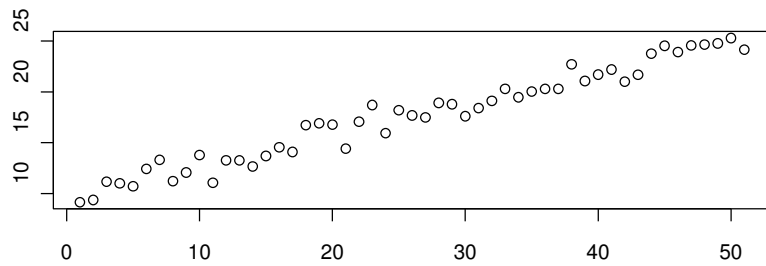


Fig.5.8-S'erieavecunetendancelinéairedependantdutemps

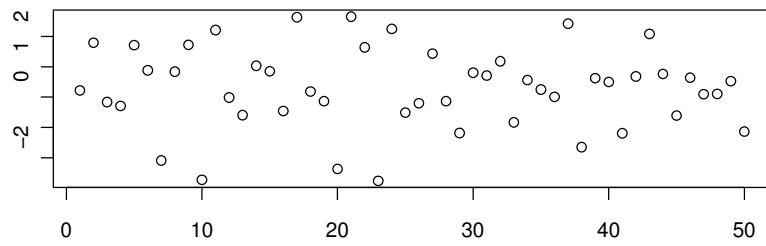


Fig.5.9-Diff'erenced'ordreundelas'erieavecunetendancelinéaire

Onpeutconstruirel'op'érateurdiff'erenced'ordredeuxen'élevant ∇ àlacarré

$\nabla_2 = \nabla \times \nabla = I - 2L + L^2$ L'op'érateurdiff'erenced'ordredeuxpermetd'enleverunetendancequadratique.Eneffet,sila

eries'écrit

$y_t = a + b \times t + c \times t^2 + E_t$, alors $\nabla_2 y_t = (I - 2L + L^2)y_t = a + b \times t + c \times t^2 + E_t - 2a - 2b \times (t-1) - 2c \times (t-1)^2 - 2E_{t-1} + a + b \times (t-2) + c \times (t-2)^2 + E_{t-1}$
'érateurs,lacomposantesaisonn`eredispara^it. **En langage R**67

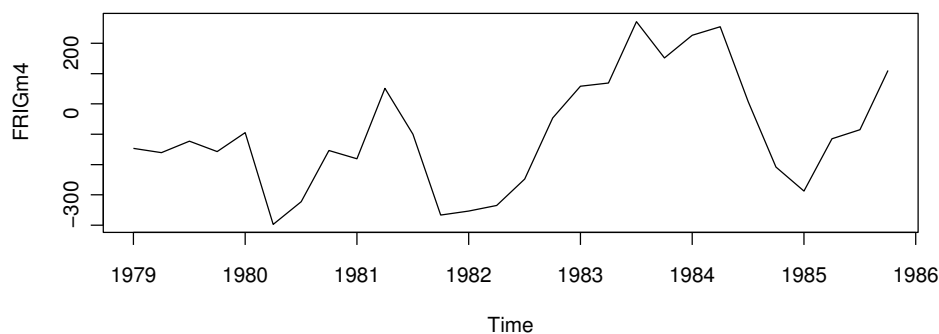


Fig.5.10-Différencé d'ordre 4 de la variable ventede réfrigérateurs

##Ventederéfrigérateursdifférencé d'ordre 4

#FRIGm4=FRIG-lag(FRIG,-4)plot(FRIGm4)

Exemple 5.7 Sion applique une différencence saisonnière d'ordre 12 sur

les données de voyageurs-kilomètres de la SNCF, on a un nombre de voyageurs-kilomètres de la SNCF qui est une variable saisonnière. On a ainsi une nouvelle variable $z_t = \nabla_{12} y_t = (I - L_{12}) y_t = y_t - y_{t-12}$. Une autre manière de faire

ce qui revient à prendre le logarithme du rapport de la variable (voir Figure 5.13). On définit ainsi une nouvelle variable $v_t = \nabla_{12} \log y_t = (I - L_{12}) \log y_t = \log y_t - \log y_{t-12}$.

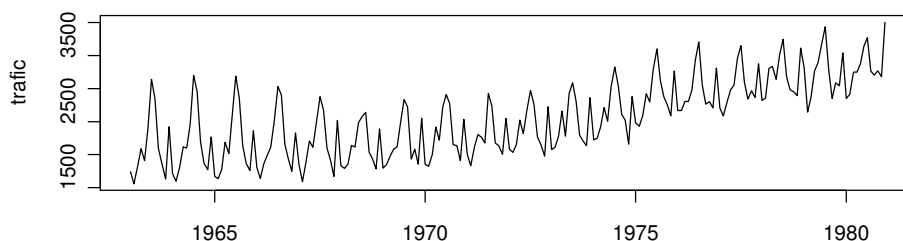
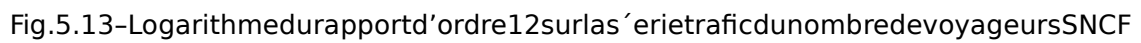


Fig.5.11-Trafic d'un nombre de voyageurs SNCF

En langage R `trafic=c(1750,1560,1820,2090,1910,2410,3140,2810,2510,2210,1910,1610,1310,1010,710,410,110,1750,1560,1820,2090,1910,2410,3140,2810,2510,2210,1910,1610,1310,1010,710,410,110)`



1834,1792,1860,2138,2115,2485,2581,2639,2038,1936,1784,2391,1798,1850,1981,2085,2120,2491,2834,
806,2976,3430,3705,3053,2764,2802,2707,3307,2706,2586,2796,2978,3053,3463,3649,3095,2839,2966,
´eairesUnfiltrelin´eaired´ordrem=p1+p2estd´efiniparFL=p2j=−p1WjL−j=w−p1Lp1+w−p1+1Lp1−1+⋯+w−1L+w0l+

où $p_1, p_2 \in \mathbb{N}$ et $w_j \in \mathbb{R}$.

5.4.2 Moyennes mobiles: définition

Une moyenne mobile d'ordre $m = p_1 + p_2 + 1$ est un filtre linéaire tel que

$$w_j = 1, \text{ pour tout } j = -p_1, \dots, p_2.$$

Beaucoup de moyennes mobiles sont des poids positifs, mais pas toutes.

Une moyenne mobile est symétrique si $p_1 = p_2 = p$, et

$$w_j = w_{-j}, \text{ pour tout } j = 1, \dots, p.$$

Une moyenne mobile symétrique est dite non-pondérée si

$$w_j = c \text{ pour tout } j = -p, \dots, p.$$

5.4.3 Moyenne mobile et composantes saisonnières

Une moyenne mobile est un outil intéressant pour lisser une série temporelle et donc pour enlever une composante saisonnière. On utilise de préférence des moyennes mobiles non pondérées d'ordre égal à la période, par exemple d'ordre 7 pour des données journalières, d'ordre 12 pour des données mensuelles. Par exemple, pour enlever la composante saisonnière due au jour de la semaine, on peut appliquer une moyenne mobile non pondérée d'ordre 7. $MM(7) = 1/7 (L_3 + L_2 + L + I + F + F_2 + F_3)$

Cette moyenne mobile accorde le même poids à chaque jour de la semaine. En effet,

$$MM(7)y_t = 1/7 (y_{t-3} + y_{t-2} + y_{t-1} + y_t + y_{t+1} + y_{t+2} + y_{t+3}).$$

Pour les composantes saisonnières d'une période paire, il est préférable d'utiliser des moyennes mobiles non pondérées. Il existe deux types de moyenne mobile centrée: - Si la période est paire et égale à m , ($m = 4$ pour les trimestrielles, $m = 12$ pour les mensuelles), la moyenne mobile est définie par $MM(m) = 1/m (L_2 + L + I + F + F_2)$. Ainsi, chaque trimestre conservé sera "perdu" aux extrémités des séries. **Exemple 5.8** La variable "réfrigérateur" est lissée grâce à une moyenne mobile d'ordre impair accordant un demi-poids aux deux extrémités. Par exemple, pour des données trimestrielles, la moyenne mobile est définie par $MM(4) = 1/8 (L_2 + 2L + 2I + 2F + F_2)$. Ainsi, chaque trimestre conservé sera "perdu" aux extrémités des séries. **Exemple 5.8** La variable "réfrigérateur" est lissée grâce à une moyenne mobile d'ordre impair accordant un demi-poids aux deux extrémités. Par exemple, pour des données trimestrielles, la moyenne mobile est définie par $MM(4) = 1/8 (L_2 + 2L + 2I + 2F + F_2)$. Ainsi, chaque trimestre conservé sera "perdu" aux extrémités des séries.

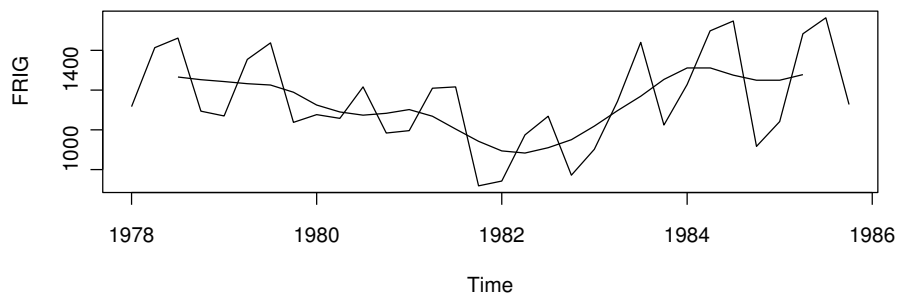


Fig.5.14-Nombre der' efrigrateur set moyennemobile d'ordre 4

En langage R `dec=decompose(FRIG) moving_average=dec$trendplot(FRIG) lines(moving_average) Unmoyen`

sonni`ere. **5.5 Moyennes mobiles particuli`eres** **5.5.1 Moyennemobile de Van Har**
 `eriodes 4 et 5 et conserve les tendances polynomiales jusqu`al'ordre 3. **5.5.3 Moyennemobile de Hender**

MoyennemobiledeHendersond'ordre $2m-3$, $m \geq 4$

$$MM_H = \sum_{j=-m-1}^{m+1} \theta_j L^j,$$

où $\theta_j =$

$$\frac{315((m-1)^2 - j^2)(m^2 - j^2)((m+1)^2 - j^2)(3m^2 - 16 - 11j^2)}{8m(m^2 - 1)(4m^2 - 1)(4m^2 - 9)(4m^2 - 25)}$$

MoyennemobiledeHendersond'ordre $2m-3=5$ ($m=4$)

$$\frac{1286}{2}(-21L_5 + 84L_4 + 160L_3 + 84F_2 - 21F_1^2)$$

MoyennemobiledeHendersond'ordre $2m-3=9$ ($m=6$)

$$\frac{12431}{2}(-99L_4 - 24L_3 - 288L_2 + 648L_1 + 805I + 648F + 288F_2 - 24F^3 - 99F^4)$$

MoyennemobiledeHendersond'ordre $2m-3=11$ ($m=7$)

$$\frac{192378}{2}(-2574L_5 - 2475L_4 + 3300L_3 + 13050L_2 + 22050L_1 + 25676I + 22050F + 13050F_2 + 3300F_3 - 2475F^4 - 2574F^5)$$

MoyennemobiledeHendersond'ordre $2m-3=15$ ($m=9$) $\frac{1193154}{2}(-2652L_7 - 4732L_6 - 2730L_5 + 4641L_4 + 16016L_3 + 37422F + 28182F_2 + 16016F_3 + 4641F_4 - 2730F_5 - 4732F_6 - 2652F^7)$

5.5.4Médianes mobiles Siles données contiennent des valeurs aberrantes ou extrêmes, on peut remplacer

un médianemobile. Par exemple la médianemobile d'ordre 5 est définie par: $Med(5)_t = \text{Médiane}(y_{t-2}, y_{t-1}, y_t, y_{t+1}, y_{t+2})$ et la tendance $S_m = 1A - 1a(Y_{am} - T_{am})$. 72

En général, on ne dispose pas d'un même nombre d'observations, pour chaque mois. On procède à un ajustement afin que la somme des composantes saisonnières soit égale à zéro :

$$S_m = S_m - \frac{1}{M} \sum_m S_m.$$

On peut ensuite procéder à la désaisonnalisation de la série par

$$Y_{am} = Y_{am} - S_m.$$

5.6.2 Méthode multiplicative

Soit une série temporelle représentée par un modèle multiplicatif du type

$$Y_{am} = T_{am} \times S_m \times E_{am}.$$

où $a = 1, \dots, A$ représente par exemple l'année et $m = 1, \dots, M$ représente par exemple le mois. La tendance est supposée connue soit par un ajustement, soit par un moyen mobile.

On isole la composante saisonnière en faisant, pour chaque mois, la moyenne des rapports entre les valeurs observées et la tendance : $S_m = \frac{1}{A} \sum_a Y_{am} / T_{am}$

$$S_m = \frac{1}{A} \sum_a \frac{Y_{am}}{T_{am}}.$$

Ensuite, on réalise un ajustement afin que la moyenne des composantes saisonnières soit égale à 1. On corrige donc les coefficients S_m par $S_m = S_m / \sum_m S_m$. La désaisonnalisation se réalise alors par une division $Y_{am} = Y_{am} / S_m$

$$Y_{am} = \frac{Y_{am}}{S_m}$$

Le Tableau 5.4 contient la variable 'vent de r'efrigerateurs', la moyenne mobile d'ordre 4, la composante saisonnière et la série désaisonnalisée au moyen de la méthode additive. Le Tableau 5.6 présente la désaisonnalisation au moyen de la méthode multiplicative. **En langage R** `deco = decompose(FRIG, type = "multiplicative") plot(deco)` 5 L'édiction à l'horizon 1, et consiste à réaliser une moyenne des valeurs passées en affectant des poids moins impor-

Tab.5.4-D'ecompositiondelavariablenFRIG,méthodeadditive

QTR	FRIG	MM	FRIG-MM	Desaison
1	1317			1442.58
2	1615			1505.13
3	1662	1466.50	195.50	1451.20
4	1295	1453.25	-158.25	1490.09
5	1271	1442.88	-171.88	1396.58
6	1555	1432.88	122.13	1445.13
7	1639	1426.50	212.50	1428.20
8	1238	1390.13	-152.13	1433.09
9	1277	1325.25	-48.25	1402.58
10	1258	1290.88	-32.88	1148.13
11	1417	1274.13	142.88	1206.20
12	1185	1283.00	-98.00	1380.09
13	1196	1302.00	-106.00	1321.58
14	1410	1268.75	141.25	1300.13
15	1417	1203.88	213.13	1206.20
16	1911	142.88-223.88		1114.09
17	943	1095.00-152.00		1068.58
18	1175	1083.25	91.75	1065.13
19	1269	1109.88	159.13	1058.20
20	973	1150.88-177.88	1168.09	
21	1102	1218.50-116.50	1227.58	
22	1344	1296.50	47.50	1234.13
23	1641	1368.88	272.13	1430.20
24	1225	1454.13-229.13	1420.09	2514
25	1429	1512.00-83.00	1554.58	2616
26	1991	1512.00	187.00	1589.13
27	1749	1475.13		

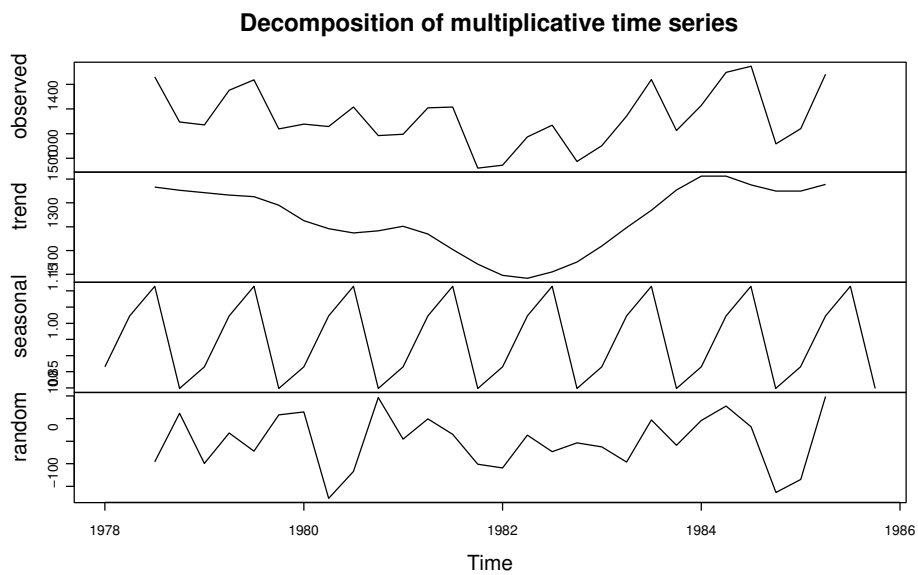
Tab.5.6-D'ecompositiondelavariablenFRIG,méthodemultiplicative

QTR	FRIG	MM	FRIG/MM	Desaison
1	1317			1453.85
2	1615			1493.76
3	1662	1466.50	1.13	1434.00
4	1295	1453.25	0.89	1516.45
5	1271	1442.88	0.88	1403.07
6	1555	1432.88	1.09	1438.26
7	1639	1426.50	1.15	1414.15
8	1238	1390.13	0.89	1449.70
9	1277	1325.25	0.96	1409.70
10	1258	1290.88	0.97	1163.56
11	1417	1274.13	1.11	1222.61
12	1185	1283.00	0.92	1387.64
13	1196	1302.00	0.92	1320.28
14	1410	1268.75	1.11	1304.15
15	1417	1203.88	1.18	1222.61
16	1911	142.88	0.80	1076.15
17	9431	1095.00	0.86	1040.99
18	1175	1083.25	1.08	1086.79
19	1269	1109.88	1.14	1094.91
20	9731	150.88	0.85	1139.39
21	1102	1218.50	0.90	1216.51
22	1344	1296.50	1.04	1243.10
23	1641	1368.88	1.20	1415.88

2412251454.130.841434.482514291512.000.951577.492616991512.001.121571.452717491475.131.191

'ee par une droite horizontale. Autrement dit, on suppose que $X_T \approx a.75$

Fig.5.15-D'ecomposition de la s rie de ventes de r frig rateurs 5.1



Le lissage exponentiel peut  tre obtenu au moyen de la m thode des moindres carr s. En annulant la d riv e par rapport   a , on obtient $\sum_{j=0}^{T-1} \beta_j (X_{T-j} - a) = 0$, ce qui donne $a = \frac{\sum_{j=0}^{T-1} \beta_j X_{T-j}}{\sum_{j=0}^{T-1} \beta_j}$. On applique alors un lissage exponentiel double pour obtenir la pr vision $\hat{X}_{T+1} = a + b$.

En annulant les dérivées partielles par rapport à a et b , on obtient

$$\sum_{j=0}^{T-1} 2\beta^j (X_{T-j} - a + bj) = 0$$

$$\sum_{j=0}^{T-1} 2\beta^j (X_{T-j} - a + bj)j = 0.$$

ce qui donne

$$\sum_{j=0}^{T-1} \beta^j X_{T-j} - a \sum_{j=0}^{T-1} \beta^j + b \sum_{j=0}^{T-1} j\beta^j = 0$$

$$\sum_{j=0}^{T-1} j\beta^j X_{T-j} - a \sum_{j=0}^{T-1} j\beta^j + b \sum_{j=0}^{T-1} j^2 \beta^j = 0.$$

Comme on a

$$\sum_{j=0}^{\infty} \beta^j = \frac{1}{1-\beta}$$

$$\sum_{j=0}^{\infty} j\beta^j = \frac{\beta}{(1-\beta)^2}$$

$$\sum_{j=0}^{\infty} j^2 \beta^j = \beta(1+\beta)(1-\beta)^{-3} \text{ on a } \sum_{j=0}^{T-1} j\beta^j X_{T-j} - a \frac{1-\beta}{1-\beta} + b \frac{\beta(1-\beta)^2}{(1-\beta)^3} = 0$$

En notant maintenant S_{1T} la série lissée $S_{1T} = (1-\beta)^{T-1} \sum_{j=0}^{T-1} \beta^j X_{T-j}$, et S_{2T} la série doublement lissée $S_{2T} = (1-\beta)^{T-1} \sum_{j=0}^{T-1} j\beta^j X_{T-j}$

Le système (5.1) peut alors s'écrire

$$\begin{cases} \frac{S_T^1}{1-\beta} - \frac{a}{1-\beta} + \frac{b\beta}{(1-\beta)^2} = 0 \\ \frac{S_T^2}{(1-\beta)^2} - \frac{S_T^1}{1-\beta} - \frac{a\beta}{(1-\beta)^2} + \frac{b\beta(1+\beta)}{(1-\beta)^3} = 0. \end{cases}$$

En résolvant ce système en a et b , on obtient finalement

$$\begin{cases} a = 2S_T^1 - S_T^2 \\ b = \frac{1-\beta}{\beta}(S_T^1 - S_T^2). \end{cases}$$

Exemple 5.10 Le tableau 5.8 rend compte du prix moyen du mazout pour 100 (achat entre 800 et 1500) en CHF pour chaque mois de 2004 à 2007 (Source: Office fédéral de la statistique, 2008).

Tab. 5.8 - Prix moyen du Mazout pour 100 (achat entre 800 et 1500)

mois/année	2004	2005	2006	2007
janvier	54.23	63.00	86.16	79.39
février	51.51	67.32	88.70	81.32
mars	55.60	75.52	88.92	82.06
avril	55.72	79.83	92.58	88.05
mai	58.71	73.22	93.65	88.24
juin	58.82	75.38	91.88	88.95
juillet	58.41	83.97	95.35	92.10
août	64.92	84.23	95.83	95.83

On obtient : $S_{11} = \hat{X}_1(1) = (1-\beta)X_1 + \beta\hat{X}_0(1) = (1-\beta)X_1 + \beta X_0(1) = (1-\beta)X_1 + \beta X_0(1)$

On obtient:

$$\begin{aligned} S_1^2 &= (1-\beta)S_1^1 + \beta S_0^2 = (1-\beta)S_1^1 + \beta S_1^1 = S_1^1 = 54.23, \\ S_2^2 &= (1-\beta)S_2^1 + \beta S_1^2 = 0.3 \times 53.414 + 0.7 \times 54.23 = 53.99, \\ S_3^2 &= (1-\beta)S_3^1 + \beta S_2^2 = 0.3 \times 54.070 + 0.7 \times 53.99 = 54.01, \end{aligned}$$

et ainsi de suite. On cherche alors

$$X_t(k) = a + bk$$

pour chaque $k=1, \dots, \hat{X}_t(1) = a + b$ avec:

$$\begin{aligned} a &= 2S_t^1 - S_t^2 \\ b &= \frac{1-\beta}{\beta} S_t^1 - S_t^2 = \frac{0.3}{0.7} S_t^1 - S_t^2 \end{aligned}$$

Le tableau 5.9 rend compte des résultats pour les années 2004 à 2007.

La figure 5.16 représente la série initiale, le lissage exponentiel simple et le lissage exponentiel double et peut être obtenue en langage R au moyen du code suivant:

```
#Lissage exponentiel double avec k=1
mazout=c(54.23,51.51,55.60,55.72,58.71,58.82,+58.41,64.92,63.95,7
+79.83,73.22,75.38,83.97,84.23,97.29,99.31,89.88,87.18,+86.16,88.70,88.92,92.58,93.65,91.88,95.35,95.
+p*liss2[i-1]#formule récursive}
#Lissage exponentiel double avec k=1
a=2*liss-liss2
b=((1-p)/p)*(liss-liss2)lis
```


Tab.5.9-Lissage exponentiel simple et double de la série temporelle Prix moyen du Mazout pour 100 litres (achat entre 800 et 1500 litres) en CHF

année	mois	X_t	S_t^1	S_t^2	a	b	$a+b$
2004	1	54.23	54.230	54.230	54.230	0.000	54.230
	2	51.51	53.414	53.985	52.843	-0.245	52.598
	3	55.60	54.070	54.011	54.129	0.025	54.154
		455.72	54.564	54.177	54.952	0.166	55.119
		558.71	55.808	54.666	56.950	0.489	57.440
		658.82	56.712	55.280	58.144	0.614	58.757
		758.41	57.221	55.862	58.580	0.582	59.163
		864.9259.531		56.963	62.099	1.101	63.199
		963.9560.857		58.131	63.582	1.168	64.750
		1072.9864.494		60.040	68.947	1.909	70.856
		1170.2566.22161.894			70.547	1.854	72.401
		1268.2466.82663.374			70.279	1.480	71.759
2005	163.0065.67864.065				67.292	0.691	67.983
		267.3266.17164.697			67.645	0.632	68.277
		375.5268.97665.98171.971				1.284	73.254
		479.8372.23267.85676.608				1.875	78.483
		573.2272.52869.25675.799				1.402	77.201
		675.3873.38470.49676.272				1.238	77.510
		783.9776.56072.31580.8051.819					82.624
		884.2378.86174.27983.4431.964					85.407
		997.2984.39077.31291.4673.033					94.501
		1099.3188.86780.77896.9533.466100.420					
2006		1189.8889.17083.29695.0442.51897.562					
		1287.1888.57384.87992.2671.58393.850					
		186.1687.84985.77089.9280.89190.819288.7088.10486.47089.7380.70090.439388.9					
		1.41380.486281.3284.03487.15680.911-1.33879.573382.0683.44186.04180.842-1.11479					

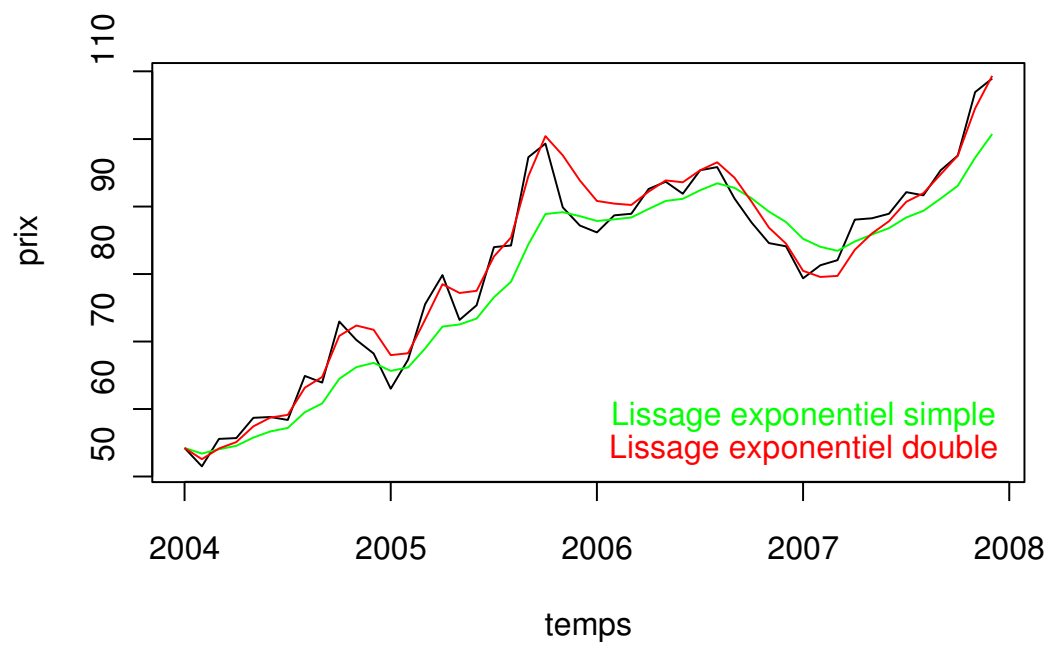


Fig.5.16–Evolution du prix du mazout en CHF (achat entre 800 et 1500), lissage exponentiel double et lissage exponentiel simple

Exercices

Exercice 5.1 Désaisonnalisiez la série suivante (c'est une série trimestrielle sur 3 années)

2417, 1605, 1221, 1826, 2367, 1569, 1176, 1742, 2804, 1399, 1063, 1755

par la méthode additive, en utilisant un moyen mobile d'ordre 4.

Exercice 5.2 En langage R utilisez la série "Ideaths" qui est une série qui se trouve dans le package de base "datasets". Lisez la documentation, puis désaisonnalisiez cette série par les méthodes additive et multiplicative. 82

Chapitre 6

Calcul des probabilités et variables

aléatoires 6.1 Probabilités 6.1.1 Événement Une expérience est dite aléatoire si son résultat ne peut être prédéterminé.

possible de cette expérience aléatoire. L'ensemble de tous les résultats possibles est noté Ω . Par exemple, si on jette deux pièces de monnaie, on peut obtenir les résultats $\Omega = \{(P,P), (F,P), (P,F), (F,F)\}$, avec F pour "face" et P pour "pile".

comme "avoir deux fois pile" ou "avoir au moins une fois pile". Formellement, un événement est un sous-ensemble de Ω . L'événement "avoir deux fois pile" est le sous-ensemble $\{(P,P)\}$. L'événement "avoir au moins une fois pile" est le sous-ensemble $\{(P,P), (F,P), (P,F)\}$.

L'ensemble Ω est appelé l'événement certain, et l'ensemble vide \emptyset est appelé l'événement impossible.

6.1.2 Opérations sur les événements Sur les événements, on peut appliquer les opérations habituelles de l'algèbre de Boole. L'événement A est "obtenir un nombre pair" et l'événement B "obtenir un multiple de 3", l'événement $A \cap B$ est l'événement "obtenir un nombre pair et multiple de 3".

Exemple 6.1 L'expérience peut consister à jeter un dé, alors

$$\Omega = \{1, 2, 3, 4, 5, 6\},$$

et un événement, A , est "obtenir un nombre pair". On a alors

$$A = \{2, 4, 6\} \text{ et } A^c = \{1, 3, 5\}.$$

6.1.3 Relations entre les événements

Événements mutuellement exclusifs

Si $A \cap B = \emptyset$ on dit que A et B sont mutuellement exclusifs, ce qui signifie que A et B ne peuvent pas se produire ensemble. **Exemple 6.2** Si on jette une pièce, l'événement "obtenir un nombre pair" et l'

événement "obtenir un nombre impair" ne peuvent pas être obtenus en même temps. Ils sont mutuellement exclusifs. Si on jette une pièce, les événements A : "obtenir un nombre pair" n'est pas mutuellement exclusif avec l'événement B : "obtenir un nombre inférieur ou égal à 3". En effet, l'intersection de A et B est non-vide et consiste en l'événement "obtenir 2". **Inclusion** Si A est inclus dans B , on écrit $A \subset B$. On dit que A implique B . **Exemple 6.3** Si on jette une pièce, les événements "obtenir un nombre pair" et "obtenir un nombre impair" sont mutuellement exclusifs.

$A = \{2\}$ et $B = \{2, 4, 6\}$. On dit que A implique B . **6.1.4 Ensemble des parties d'un ensemble et système d'événements** Soit Ω un ensemble fini. Soit $\mathcal{P}(\Omega)$ l'ensemble des parties de Ω . Soit \mathcal{A} un système d'événements. Soit $P(\cdot)$ une application de \mathcal{A} dans $[0, 1]$, telle que: $-P(\Omega) = 1$, $-$ Pour tout ensemble d'événements A_1, A_2, \dots deux à deux incompatibles, on a: $P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$

Tab.6.1-Syst` emecomplet d'evenements

A_1	1111111	A_i	1111111	A_n
-------	---------	-------	---------	-------

Apartir des axiomes, on peut d'eduire les proprietés suivantes:

Propri´ et´ e6.1 $Pr(\emptyset)=0$. *D´emonstration* Comme \emptyset est d'intersection vide avec \emptyset , on a que $Pr(\emptyset \cup \emptyset)=Pr(\emptyset)+Pr(\emptyset)$

Donc, $Pr(\emptyset)=2Pr(\emptyset)$, ce qui implique que $Pr(\emptyset)=0$. *

Propri´ et´ e6.2 $Pr(A^c)=1-Pr(A)$. *D´emonstration* On sait que $A \cup A^c=\Omega$ et $A \cap A^c=\emptyset$. Ainsi, on a que $Pr(\Omega)=Pr(A \cup A^c)=Pr(A)+Pr(A^c)$

Propri´ et´ e6.3 $Pr(A) \leq Pr(B)$ si $A \subset B$. *D´emonstration* Comme $A \subset B$, on a $B=(B \cap A) \cup A$. Mais on a que $(B \cap A) \cap A=\emptyset$. Ainsi,

Démonstration

On a $A \cup B = A \cup (B \cap A^c)$,

avec $A \cap (B \cap A^c) = \emptyset$.

Ainsi $\Pr(A \cup B) = \Pr(A) + \Pr(B \cap A^c)$.

Il reste à montrer que $\Pr(B) = \Pr(B \cap A^c) + \Pr(B \cap A)$

Mais $B = (B \cap A^c) \cup (B \cap A)$

avec $(B \cap A^c) \cap (B \cap A) = \emptyset$

Donc $\Pr(B) = \Pr(B \cap A^c) + \Pr(B \cap A)$

*

Propriété 6.5 $\Pr_{i=1}^n A_i \leq \Pr(A_i)$ *Démonstration* Notons respectivement $B_1 = A_1, B_2 = (A_2 | A_1), B_3 = (A_3 | (A_1 \cup A_2))$,

$B_4 = (A_4 | (A_1 \cup A_2 \cup A_3)), \dots, B_n = (A_n | (A_1 \cup A_2 \cup A_3 \cup \dots \cup A_{n-1}))$.

Comme $\Pr_{i=1}^n A_i = \Pr_{i=1}^n B_i$, et que $B_i \cap B_j = \emptyset$ pour tout $j \neq i$, alors $\Pr_{i=1}^n B_i = \Pr(B_i)$. De plus, comme, pour tout $i, B_i \subset A_i$, on a $\Pr(B_i) \leq \Pr(A_i)$.

6.1.6 Probabilités conditionnelles et indépendance

Définition 6.3 Soient deux événements A et B , si $\Pr(B) > 0$, alors

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)}.$$

Exemple 6.5 Soit une urne, et qu'on considère les deux événements suivants:

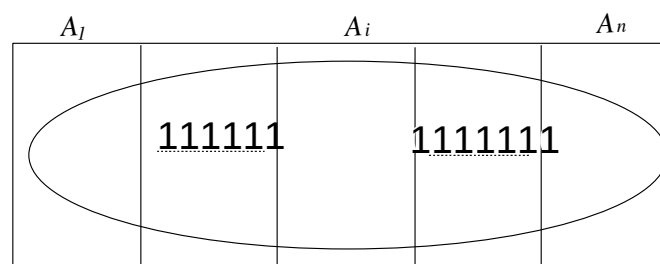
- A la probabilité d'avoir un nombre pair
- B la probabilité d'avoir un nombre supérieur ou égal à 4.

On a donc $\Pr(A) = \Pr(\{2, 4, 6\}) = 12$, $\Pr(B) = \Pr(\{4, 5, 6\}) = 36 = 12$, $\Pr(A \cap B) = \Pr(\{4, 6\}) = 26 = 13$, $\Pr(A|B) = \Pr(A \cap B) / \Pr(B)$



$\Pr(A|B) = \Pr(A)$. On peut montrer facilement que si A et B sont indépendants, alors $\Pr(A \cap B) = \Pr(A)\Pr(B)$. **6.1.7 Théorème**

$\Pr(B) = \sum_{i=1}^n \Pr(A_i) \Pr(B|A_i)$. Tab. 6.2 – Illustration du théorème des probabilités totales



En effet, $\sum_{i=1}^n \Pr(A_i) \Pr(B|A_i) = \sum_{i=1}^n \Pr(B \cap A_i)$. Comme les événements $A_i \cap B$ sont mutuellement exclusifs, $\sum_{i=1}^n \Pr(B \cap A_i) =$

Théorème 6.2 (de Bayes) Soit A_1, \dots, A_n un système complet d'événements, alors

$$\Pr(A_i|B) = \frac{\Pr(A_i)\Pr(B|A_i)}{\sum_{j=1}^n \Pr(A_j)\Pr(B|A_j)}.$$

En effet, par le théorème des probabilités totales,

$$\frac{\Pr(A_i)\Pr(B|A_i)}{\sum_{j=1}^n \Pr(A_j)\Pr(B|A_j)} = \frac{\Pr(B \cap A_i)}{\Pr(B)} = \Pr(A_i|B).$$

Exemple 6.6 Supposons qu'une population d'adultes soit composée de 30% de fumeurs (A_1) et de 70% de non-fumeur (A_2). Notons B l'événement "mourir d'un cancer du poumon". Supposons en outre que la probabilité de mourir d'un cancer du poumon est égale à $\Pr(B|A_1) = 20\%$ si l'on est fumeur et de $\Pr(B|A_2) = 1\%$ si l'on est non-fumeur. Le théorème de Bayes permet de calculer les probabilités a priori, c'est-à-dire la probabilité d'avoir été fumeur si on est mort d'un cancer du poumon. En effet, cette probabilité est notée $\Pr(A_1|B)$ et peut être calculée par $\Pr(A_1|B) = \frac{\Pr(A_1)\Pr(B|A_1)}{\Pr(A_1)\Pr(B|A_1) + \Pr(A_2)\Pr(B|A_2)} =$

$$\frac{0.3 \times 0.2}{0.3 \times 0.2 + 0.7 \times 0.01} = \frac{0.06}{0.06 + 0.007} \approx 0.896.$$

La probabilité de ne pas avoir été non-fumeur si on est mort d'un cancer du poumon vaut quant à elle:

$$\Pr(A_2|B) = \frac{\Pr(A_2)\Pr(B|A_2)}{\Pr(A_1)\Pr(B|A_1) + \Pr(A_2)\Pr(B|A_2)} = \frac{0.7 \times 0.01}{0.3 \times 0.2 + 0.7 \times 0.01} = \frac{0.007}{0.06 + 0.007} \approx 0.104.$$

6.2 Analyse combinatoire 6.2.1 Introduction L'analyse combinatoire est l'étude mathématique

natoire est un outil utilisé dans le calcul des probabilités. **6.2.2 Permutations (sans répétition)** Une

l'ensemble $\{1, 2, 3\}$. Il existe 6 manières d'ordonner ces trois chiffres: $\{1, 2, 3\}, \{1, 3, 2\}, \{2, 1, 3\}, \{2, 3, 1\}, \{3, 1, 2\}, \{3, 2, 1\}$.

6.2.3 Permutations avec répétition

On peut également se poser la question d'un nombre de manières d'arranger des objets qui ne sont pas tous distincts. Supposons que nous ayons 2 boules rouges (notées R) et 3 boules blanches (notées B). Il existe 10 permutations possibles qui sont :

$\{R, R, B, B, B\}, \{R, B, R, B, B\}, \{R, B, B, R, B\}, \{R, B, B, B, R\}, \{B, R, R, B, B\},$

$\{B, R, B, R, B\}, \{B, R, B, B, R\}, \{B, B, R, R, B\}, \{B, B, R, B, R\}, \{B, B, B, R, R\}.$

Sil'on dispose de n objets appartenant à deux groupes détaillés de n_1 et n_2 , le nombre de permutations avec répétition est $n!$

$$\frac{n!}{n_1!n_2!}.$$

Par exemple si l'on a 3 boules blanches et 2 boules rouges, on obtient

$$\frac{n!}{n_1!n_2!} = \frac{5!}{2!3!} = \frac{120}{2 \times 6} = 10.$$

Sil'on dispose de n objets appartenant à p groupes détaillés de n_1, n_2, \dots, n_p , le nombre de permutations avec répétition est $n!$

6.2.4 Arrangements (sans répétition)

Soit n objets distincts. On appelle un arrangement une manière de les ranger dans des boîtes numérotées de 1 à k . Dans la première boîte, on peut mettre chacun des n objets. Dans la deuxième, on peut mettre chacun des $n-1$ objets restants, dans la troisième, on peut mettre chacun des $n-2$ objets restants et ainsi de suite. Le nombre d'arrangements possibles est donc égal à : $A_{kn} = n \times (n-1) \times (n-2) \times \dots \times (n-k+1) = \frac{n!}{(n-k)!}$

Le nombre de combinaisons de k objets parmi n est le nombre de sous-ensembles de taille k dans un ensemble de taille n . Soit l'ensemble $\{1, 2, 3, 4, 5\}$. Il existe 10 sous-ensembles de taille 3 qui sont : $\{1, 2, 3\}, \{1, 2, 4\}, \{1, 2, 5\}, \{1, 3, 4\}, \{1, 3, 5\}, \{1, 4, 5\}, \{2, 3, 4\}, \{2, 3, 5\}, \{2, 4, 5\}, \{3, 4, 5\}$. Le nombre de combinaisons de k objets parmi n est noté C_{kn} et est égal à : $C_{kn} = \frac{n!}{k!(n-k)!}$

6.3 Variables aléatoires

6.3.1 Définition

La notion de variable aléatoire formalise l'association d'une valeur au résultat d'une expérience aléatoire.

Définition 6.5 Une variable aléatoire X est une application de l'ensemble fondamental Ω dans \mathbb{R} .

Exemple 6.7 On considère une expérience aléatoire consistant à lancer deux pièces de monnaie. L'ensemble des résultats possibles est $\Omega = \{(F,F), (F,P), (P,F), (P,P)\}$.

Chacun des éléments de Ω a une probabilité $1/4$. Une variable aléatoire va associer une valeur à chacun des éléments de Ω . Considérons la variable aléatoire représentant le nombre de "Faces" obtenus:

$X = \begin{cases} 0 & \text{avec une probabilité } 1/4 \\ 1 & \text{avec une probabilité } 1/2 \\ 2 & \text{avec une probabilité } 1/4. \end{cases}$

C'est une variable aléatoire discrète dont la distribution de probabilité est présentée en Figure 6.1.

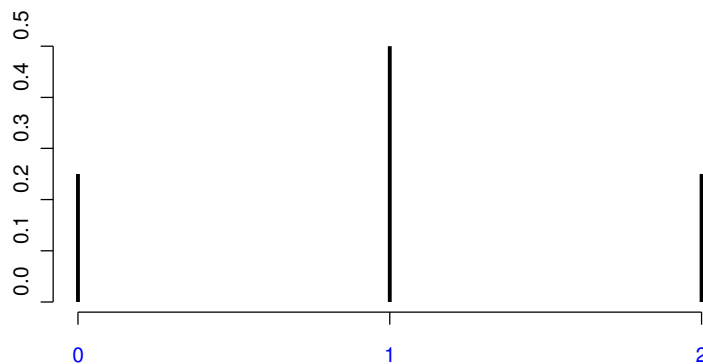


Fig. 6.1 - Distribution de "faces" obtenus.

6.4 Variables aléatoires discrètes

6.4.1 Définition

$(X - \mu)^2 = \sum_{x \in \mathbb{Z}} p_X(x) x^2 - \mu^2$. On peut aussi calculer les moments et tous les autres paramètres. 90

6.4.2 Variable indicatrice ou bernoullienne

La variable indicatrice X de paramètre $p \in [0, 1]$ a la distribution de probabilité suivante :

$$X = \begin{cases} 1 & \text{avec une probabilité } p \\ 0 & \text{avec une probabilité } 1-p. \end{cases}$$

L'espérance vaut $\mu = E(X) = 0 \times (1-p) + 1 \times p = p$,

et la variance vaut $\sigma^2 = \text{var}(X) = E(X-p)^2$

$$= (1-p)(0-p)^2 + p(1-p)^2 = p(1-p).$$

Exemple 6.8 On tire au hasard une boule dans une urne contenant 18 boules rouges et 12 boules blanches. Si

X vaut 1 si la boule est rouge et 0 sinon, alors X a une loi bernoullienne de paramètre $p = 18/(18+12) = 0.6$.

6.4.3 Variable binomiale

La variable aléatoire binomiale de paramètres n et p correspond à l'expérience suivante : On renouvelle n fois de manière indépendante une épreuve de Bernoulli de paramètre p , où p est la probabilité de succès pour une expérience élémentaire. Ensuite, on note X le nombre de succès obtenus. Le nombre de succès est une variable aléatoire prenant des valeurs entières de 0 à n et ayant une distribution binomiale.

Une variable X suit une loi binomiale de paramètres $0 < p < 1$ et d'exposant n , si

$$\Pr(X=x) = n \times p^x \times q^{n-x}, \quad x=0, 1, \dots, n-1, n,$$

où $q=1-p$, et $n! = n \times (n-1) \times \dots \times 1$. Demandez synthétique, si X a une distribution binomiale, on note : $X \sim B(n, p)$. **Rappel**

du binôme de Newton $(p+q)^n$. $(p+q)^0 = 1$, $(p+q)^1 = p+q$, $(p+q)^2 = p^2 + 2pq + q^2$, $(p+q)^3 = p^3 + 3p^2q + 3pq^2 + q^3 = 1$

L'espérance calculée de la manière suivante:

$$\begin{aligned}
 E(X) &= \sum_{x=0}^n x \Pr(X=x) \\
 &= \sum_{x=0}^n x \binom{n}{x} p^x q^{n-x} \\
 &= \sum_{x=1}^n x \binom{n}{x} p^x q^{n-x} \quad (\text{on peut enlever le terme } x=0) \\
 &= \sum_{x=1}^n \binom{n-1}{x-1} p^x q^{n-x} \\
 &= np \sum_{x=1}^n \binom{n-1}{x-1} p^{x-1} q^{(n-1)-(x-1)} \\
 &= np \sum_{z=0}^{n-1} \binom{n-1}{z} p^z q^{(n-1)-z} \quad (\text{en posant } z=x-1)
 \end{aligned}$$

$= np(p+q)^{n-1} = np$. La variance est donnée (sans démonstration) par $\text{var}(X) = npq$. **Exemple 6.9** On tire au hasard

18 boules rouges et 12 boules blanches. Si X est le nombre de boules rouges obtenues, alors X a un loi binomiale de paramètre $p = 18/(18+12) = 0.6$, et d'exposant $n = 5$. Donc, $\Pr(X=x) = \binom{5}{x} 0.6^x 0.4^{5-x}$, $x=0,1,\dots,4,5$, ce qui est dans la Figure 6.2. **Exemple 6.10** Supposons que, dans une population d'électeurs, 60% des électeurs s'apprêtent

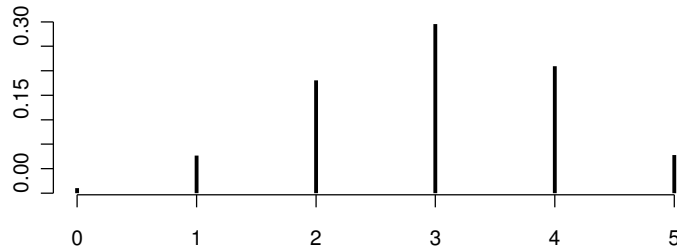


Fig.6.2-Distribution d'une variable aléatoire binomiale avec $n=5$ et $p=0.6$.

6.4.4 Variable de Poisson

La variable X suit une loi de Poisson, de paramètre $\lambda \in \mathbb{R}^+$ si

$$\Pr(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}, x=0, 1, 2, 3, \dots$$

On note alors $X \sim P(\lambda)$. La somme des probabilités est bien égale à 1, en effet

$$\sum_{x=0}^{\infty} \Pr(X=x) = \sum_{x=0}^{\infty} e^{-\lambda} \frac{\lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1.$$

L'espérance et la variance d'une loi de Poisson sont égales au paramètre λ . En effet

$$E(X) = \sum_{x=0}^{\infty} x \Pr(X=x) = \sum_{x=0}^{\infty} x e^{-\lambda} \frac{\lambda^x}{x!} = e^{-\lambda} \sum_{x=1}^{\infty} x \lambda^{x-1} (x-1)! = e^{-\lambda} \sum_{z=0}^{\infty} \lambda^z z! = \lambda e^{-\lambda} \sum_{z=0}^{\infty} \lambda^z z! = \lambda e^{-\lambda} e^{\lambda} = \lambda.$$

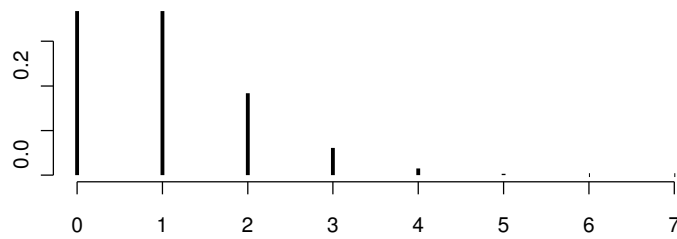


Fig.6.3-Distribution d'une variable de Poisson avec $\lambda=1$.

```
lwd=3,xlab="",ylab="",main="",frame=FALSE)
#PoissonP(1)plot(dpois(0:7,1),type="h",lwd=3,xlab="",ylab="",main="",frame=FALSE)
```

6.5 Variable aléatoire continue 6.5.1 Définition, espérance et variance

La probabilité qu'une variable aléatoire continue soit inférieure à une valeur donnée se définit par la fonction de répartition. $\Pr(X \leq x) = F(x)$. La fonction de répartition d'une variable aléatoire continue est toujours continue et croissante. La probabilité qu'une variable aléatoire continue soit comprise entre deux valeurs est donnée par la différence des fonctions de répartition. $\Pr(a < X < b) = F(b) - F(a)$. Dans la Figure 6.4, la probabilité $\Pr[X \leq a]$ est l'aire sous la courbe de la fonction de densité de probabilité $f(x)$ pour $x \leq a$.

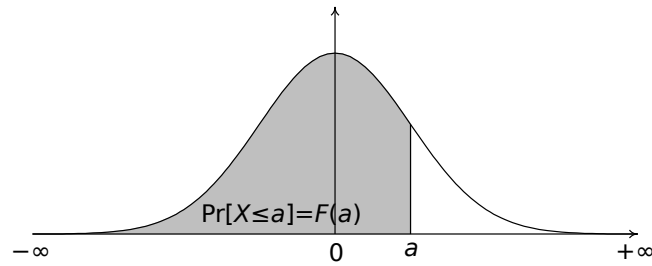


Fig.6.4-Probabilité que la variable aléatoire soit inférieure à a

Si la variable aléatoire est continue, la probabilité qu'elle prenne exactement une valeur quelconque est nulle:

$$\Pr[X=a]=0.$$

L'espérance d'une variable aléatoire continue est définie par:

$$E(X) = \int_{-\infty}^{+\infty} x f(x) dx,$$

et la variance $\text{var}(X) = \int_{-\infty}^{+\infty} (x-\mu)^2 f(x) dx$.

6.5.2 Variable uniforme Une variable aléatoire X est dite uniforme dans un intervalle $[a, b]$, (avec $a < b$) si sa répartition est:

$F(x) = \begin{cases} 0 & \text{si } x < a \\ (x-a)/(b-a) & \text{si } a \leq x \leq b \\ 1 & \text{si } x > b \end{cases}$. Sa densité est alors $f(x) = \begin{cases} 0 & \text{si } x < a \\ 1/(b-a) & \text{si } a \leq x \leq b \\ 0 & \text{si } x > b \end{cases}$. On peut

*

Résultat 6.2

$$\sigma^2 = \text{var}(X) = \frac{(b-a)^2}{12}.$$

Démonstration Demandez en exercice, une variance peut toujours s'écrire comme un moment d'ordre 2 moins le carré de la moyenne. En effet,

$$\begin{aligned} \sigma^2 &= \text{var}(X) \\ &= \int_a^b (x - \mu)^2 f(x) dx \\ &= \int_a^b (x^2 + \mu^2 - 2x\mu) f(x) dx \\ &= \int_a^b x^2 f(x) dx + \int_a^b \mu^2 f(x) dx - 2\mu \int_a^b x f(x) dx \\ &= \int_a^b x^2 f(x) dx + \mu^2 - 2\mu \int_a^b x f(x) dx \end{aligned}$$

$= \int_a^b x^2 f(x) dx - \mu^2$. On calcule ensuite un moment d'ordre 2 : $\int_a^b x^2 f(x) dx = \int_a^b x^2 (1/b - a/x) dx = 1/b - a \int_a^b x dx = 1/b - a \frac{x^2}{2} \Big|_a^b = 1/b - a \frac{b^2 - a^2}{2}$

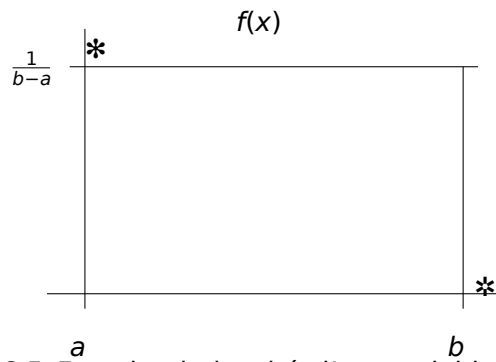


Fig.6.5-Fonction de densité d'une variable uniforme

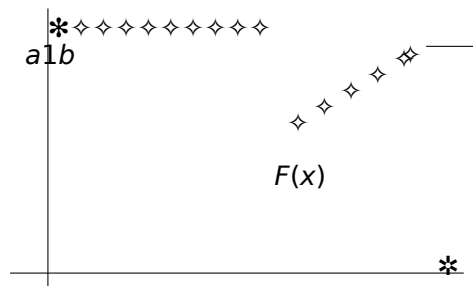


Fig.6.6-Fonction de répartition d'une variable uniforme

6.5.3 Variable normale Une variable aléatoire X est dite normale si sa densité $f_{\mu, \sigma^2}(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

$$f_{\mu, \sigma^2}(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

où $\mu \in \mathbb{R}$ et $\sigma \in \mathbb{R}_+$ sont les paramètres de la distribution. Le paramètre μ est appelé la moyenne et le paramètre σ l'écart-type de la variable normale. $\mu - \infty < \mu + \sigma$

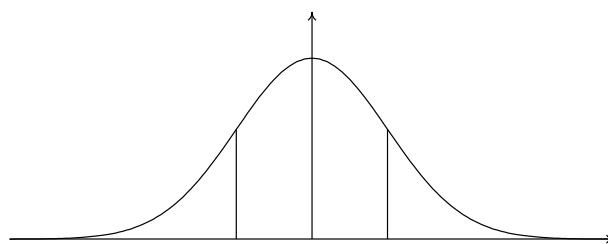


Fig.6.7-Fonction de densité d'une variable normale

La fonction de répartition vaut

$$F_{\mu, \sigma^2}(x) = \int_{-\infty}^x \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{u-\mu}{\sigma}^2\right) du.$$

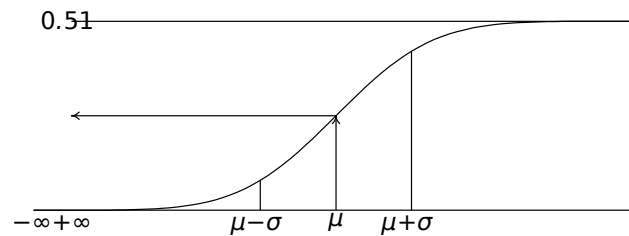
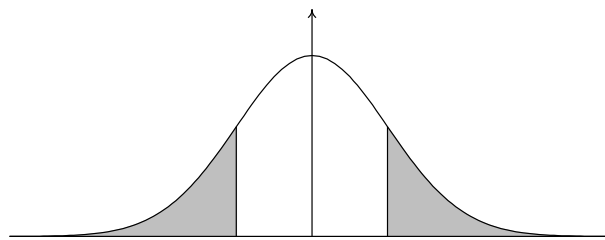


Fig.6.8-Fonction de répartition d'une variable normale

6.5.4 Variable normale centrée et réduite La variable aléatoire normale centrée et réduite est une variable d'espérance nulle $\mu=0$ et de variance $\sigma^2=1$. Sa fonction de densité vaut $f_{0,1}(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$ et sa fonction de répartition vaut $\Phi(x) = F_{0,1}(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du$.



En posant

$$z = \frac{u - \mu}{\sigma},$$

on obtient $u = z\sigma + \mu$, et donc $du = \sigma dz$. Donc,

$$F_{\mu, \sigma^2}(x) = \int_{-\infty}^{\frac{x-\mu}{\sigma}} \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \sigma dz = \Phi\left(\frac{x-\mu}{\sigma}\right).$$

Lestables de la variable normale ne sont données que pour la normale centrée réduite. Lestables ne donnent $\Phi(x)$ que pour les valeurs positives de x , car les valeurs négatives peuvent être retrouvées par la relation de symétrie. **6.5.5 Distribution exponentielle** *

Une variable aléatoire X a une distribution exponentielle si sa fonction de densité est donnée par:

$$f(x) = \begin{cases} \lambda \exp(-\lambda x), & \text{si } x > 0 \\ 0 & \text{sinon} \end{cases}$$

Le paramètre λ est positif. Quand $x > 0$, sa fonction de répartition vaut:

$$F(x) = \int_0^x f(u) du = \int_0^x \lambda e^{-\lambda u} du = -e^{-\lambda u} \Big|_0^x = 1 - e^{-\lambda x}.$$

On peut alors calculer la moyenne: **Résultat 6.4** $E(X) = \frac{1}{\lambda}$ **Démonstration** $E(X) = \int_0^\infty x f(x) dx = \int_0^\infty x \lambda e^{-\lambda x} dx = -\frac{1}{\lambda} + x e^{-\lambda x} \Big|_0^\infty = \frac{1}{\lambda}$

Il est également possible de montrer que la variance vaut: $\text{var}(X) = \frac{1}{\lambda^2}$. **6.6 Distribution bivariée** *

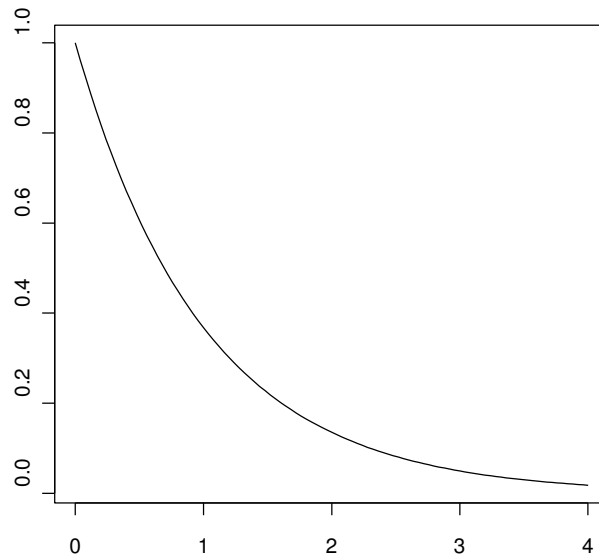


Fig.6.10-Fonction de densité d'une variable exponentielle avec $\lambda = 1$.

Avec les distributions marginales, on peut définir les moyennes marginales, et les variances marginales:

$$\mu_X = \int_{-\infty}^{\infty} x f_X(x) dx, \text{ et } \mu_Y = \int_{-\infty}^{\infty} y f_Y(y) dy,$$

$$\sigma_{2X} = \int_{-\infty}^{\infty} (x - \mu_X)^2 f_X(x) dx, \text{ et } \sigma_{2Y} = \int_{-\infty}^{\infty} (y - \mu_Y)^2 f_Y(y) dy.$$

On appelle densités conditionnelles, les fonctions $f(x|y) = f(x, y) / f_Y(y)$ et $f(y|x) = f(x, y) / f_X(x)$. Avec les distributions conditionnelles:

$\mu_X(y) = \int_{-\infty}^{\infty} x f(x|y) dx$, et $\mu_Y(x) = \int_{-\infty}^{\infty} y f(y|x) dy$, $\sigma_{2X}(y) = \int_{-\infty}^{\infty} \{x - \mu_X(y)\}^2 f(x|y) dx$, et $\sigma_{2Y}(x) = \int_{-\infty}^{\infty} \{y - \mu_Y(x)\}^2 f(y|x) dy$.
 Si on joint les deux variables, alors X et Y sont indépendantes si $f_{XY}(x, y) = f_X(x) f_Y(y)$, $x, y \in \mathbb{R}$. 100

6.7 Propriétés des espérances et des variances

De manière générale, pour des variables aléatoires X et Y , et avec a et b constants, on a les résultats suivants. **Résultat 6.5** $E(a + bX) = a + bE(X)$

Démonstration $E(a + bX) = \int_{\mathbb{R}} (a + bx)f(x)dx = a$

$$\int_{\mathbb{R}} f(x)dx + b \int_{\mathbb{R}} xf(x)dx = a + bE(X).$$

*

Résultat 6.6 $E(aY + bX) = aE(Y) + bE(X)$

Démonstration $E(aY + bX) = \int_{\mathbb{R}^2} (ay + bx)f(x, y)dxdy$

$$= a \int_{\mathbb{R}^2} yf(x, y)dxdy + b \int_{\mathbb{R}^2} xf(x, y)dxdy$$

$$= a \int_{\mathbb{R}} y \left(\int_{\mathbb{R}} f(x, y)dx \right) dy + b \int_{\mathbb{R}} x \left(\int_{\mathbb{R}} f(x, y)dy \right) dx$$

$$= a \int_{\mathbb{R}} yf(y)dy + b \int_{\mathbb{R}} xf(x)dx = aE(Y) + bE(X) *$$

Résultat 6.7 $\text{var}(a + bX) = b^2 \text{var}(X)$. **Démonstration** $\text{var}(a + bX) = \int_{\mathbb{R}} [a + bx - E(a + bX)]^2 f(x)dx = \int_{\mathbb{R}} [a + bx - (a + bE(X))]$

Démonstration

$$\begin{aligned}
 \text{var}(X+Y) &= \int_{\mathbb{R}} \int_{\mathbb{R}} [x+y-E(X+Y)]^2 f(x,y) dx dy \\
 &= \int_{\mathbb{R}} \int_{\mathbb{R}} [x-E(X)+y-E(Y)]^2 f(x,y) dx dy \\
 &= \int_{\mathbb{R}} \int_{\mathbb{R}} [x-E(X)]^2 + [y-E(Y)]^2 + 2[x-E(X)][y-E(Y)] f(x,y) dx dy \\
 &= \text{var}(X) + \text{var}(Y) + 2\text{cov}(X,Y)
 \end{aligned}$$

*

Résultat 6.9 De plus, si X et Y sont indépendantes, on a $f(x,y) = f(x)f(y)$ pour tout x, y

$$E(XY) = E(X)E(Y).$$

Démonstration $E(XY) = \int_{\mathbb{R}} \int_{\mathbb{R}} xy f(x)f(y) dx dy$

$$\begin{aligned}
 &= \int_{\mathbb{R}} x f(x) dx \int_{\mathbb{R}} y f(y) dy \\
 &= E(X)E(Y). *
 \end{aligned}$$

Enfin, si X et Y sont indépendantes, on a $\text{cov}(X,Y) = 0$, et donc $\text{var}(X+Y) = \text{var}(X) + \text{var}(Y)$. Enfin, il est possible de calculer

et identiquement distribuées. **Théorème 6.3** Soit X_1, \dots, X_n une suite de variables aléatoires indépendantes et identiquement distribuées.

et dont la moyenne μ et la variance σ^2 existent et sont finies, alors si $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$, on a $E(\bar{X}) = \mu$, et $\text{var}(\bar{X}) = \frac{\sigma^2}{n}$. **Démonstration**

—

—

—

—

—

—

—

—

—

6.8 Autres variables aléatoires

6.8.1 Variable khi-carrée

Soit une suite de variables aléatoires indépendantes, normales, centrées et réduites, X_1, \dots, X_p , (c'est-à-dire de moyenne nulle et de variance égale à 1), alors la variable aléatoire

$$\chi_p^2 = \sum_{i=1}^p X_i^2,$$

est appelée variable aléatoire khi-carrée à p degrés de liberté. Il est possible de montrer que $E(\chi_p^2) = p$,

et que $\text{var}(\chi_p^2) = 2p$.

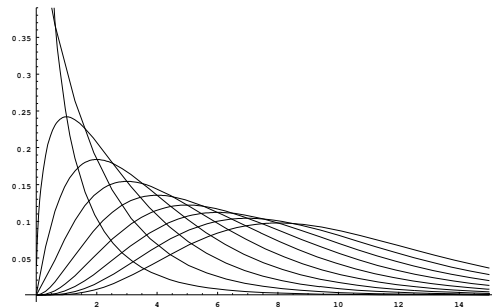


Fig.6.11-Densité d'une variable de chi-carrée avec $p=1, 2, \dots, 10$

6.8.2 Variable de Student

Soit une variable aléatoire normale centrée et réduite, X , et une variable aléatoire indépendante de X , alors la variable aléatoire $t_p = X / \sqrt{\chi_p^2 / p}$ est appelée variable aléatoire de Student à p degrés de liberté.

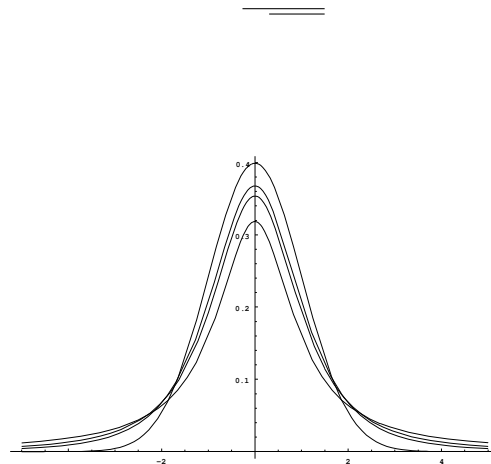


Fig.6.12-Densité de variables de Student avec $p=1, 2$ et 3 et d'une variable normale

6.8.3 Variable de Fisher

Soient deux variables aléatoires khi-carrées indépendantes χ_p^2 et χ_q^2 , respectivement à p et q degrés de liberté, alors la variable aléatoire

$$F_{p,q} = \frac{\chi_p^2/p}{\chi_q^2/q}$$

est appelée variable de Fisher à p et q degrés de liberté.

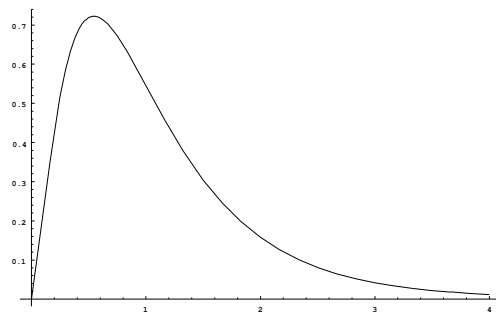


Fig.6.13–Densité d'une variable de Fisher

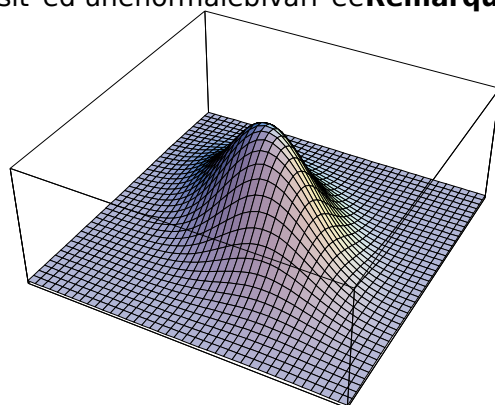
Remarque 6.1 Il est facile de montrer que le carré d'une variable de Student à q degrés de liberté est une variable de Fisher à 1 et q degrés de liberté.

6.8.4 Variable normale multivariée

Soit \mathbf{x} un vecteur de variable aléatoire de dimension p à densité de probabilité $f(\mathbf{x})$ et de moyenne $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$ et de matrice variance-covariance $\boldsymbol{\Sigma}$ (on suppose pour simplifier que $\boldsymbol{\Sigma}$ est de plein rang), si sa fonction de densité est donnée par $f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$, (6.2)

pour tout $\mathbf{x} \in \mathbb{R}^p$.

Fig.6.14–Densité d'une normale bivariée



Remarque 6.2 Si $p=1$, on retrouve l'expression (6.1).

Un cas particulier est important: supposons qu'une matrice variance-covariance peut s'écrire $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$, ce qui signifie que toutes les composantes du vecteur \mathbf{X} sont non-corrélées. Dans ce cas,

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \frac{1}{(2\pi)^{p/2} (\prod_{j=1}^p \sigma_j^2)^{1/2}} \exp - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \frac{1}{(2\pi)^{p/2} (\prod_{j=1}^p \sigma_j)} \exp \left[- \sum_{j=1}^p \frac{(x_j - \mu_j)^2}{2\sigma_j^2} \right] \\ &= \frac{1}{(2\pi)^{p/2} (\prod_{j=1}^p \sigma_j)} \exp - \sum_{j=1}^p \frac{(x_j - \mu_j)^2}{2\sigma_j^2} \\ &= \prod_{j=1}^p \frac{1}{(2\pi)^{1/2} \sigma_j} \exp - \frac{(x_j - \mu_j)^2}{2\sigma_j^2} \end{aligned}$$

$$= \prod_{j=1}^p f_{X_j}(x_j), \text{ où } f_{X_j}(x_j) = \frac{1}{(2\pi\sigma_j^2)^{1/2}} \exp - \frac{(x_j - \mu_j)^2}{2\sigma_j^2}$$

$$\frac{1}{(2\pi\sigma_j^2)^{1/2}} \exp - \frac{(x_j - \mu_j)^2}{2\sigma_j^2},$$

est la densité de la variable X_j . On constate que s'il y a absence de corrélation entre les variables normales, alors la densité du vecteur normal peut s'écrire comme un produit de densités univariées. Dans le cas multinormal (et seulement dans ce cas), l'absence de corrélation implique donc l'indépendance des variables.

De manière générale, si \mathbf{X} est un vecteur de variables aléatoires de moyenne $\boldsymbol{\mu}$ et de matrice variance-covariance Σ , et si \mathbf{A} est une matrice $q \times p$ de constantes, alors $E(\mathbf{AX}) = \mathbf{A}E(\mathbf{X}) = \mathbf{A}\boldsymbol{\mu}$, et $\text{var}(\mathbf{AX}) = \mathbf{A}\text{var}(\mathbf{X})\mathbf{A}^T = \mathbf{A}\Sigma\mathbf{A}^T$. Dans

pendant, la matrice variance-covariance n'est pas nécessairement de plein rang. Donc, si \mathbf{X} est un vecteur multinormal

Exercice 6.2 Déterminez les valeurs j de la variable normale centrée réduite Z telles que:

1. $\Pr[Z \leq j] = 0,9332$;
2. $\Pr[-j \leq Z \leq j] = 0,3438$;
3. $\Pr[Z \leq j] = 0,0125$;
4. $\Pr[Z \geq j] = 0,0125$;
5. $\Pr[j \leq Z \leq 3] = 0,7907$.

Exercice 6.3 Soit une variable aléatoire $X \sim N(53; \sigma^2 = 100)$ représentant le résultat d'un examen pour un étudiant d'une section. Déterminez la probabilité pour que le résultat soit compris entre 33,4 et 72,6.

Exercice 6.4 Soit une variable aléatoire $X \sim N(50; \sigma^2 = 100)$. Déterminez le premier quartile de cette distribution. **Exercice 6.5** En supposant que les tailles en cm des étudiants d'un pays admettent la distribution normale

$N(172; \sigma^2 = 9)$. On demande de déterminer le pourcentage théorique:

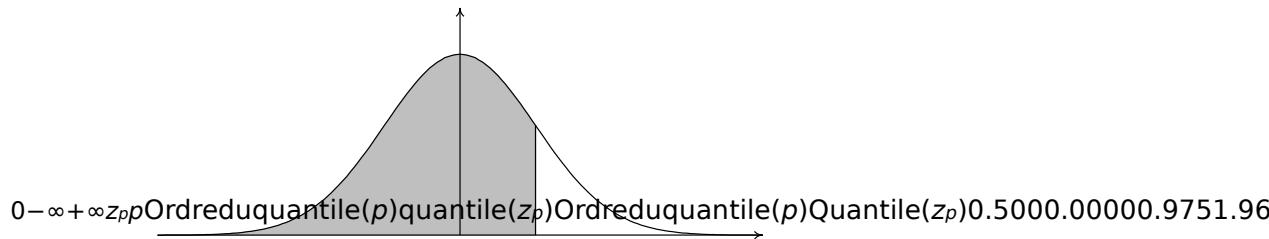
a) d'étudiants mesurant au moins 180 cm. b) d'étudiants dont la taille est comprise entre 168 et 180. **Exercice 6.6** Soit

la vitesse de toutes les automobiles pendant une journée. En supposant que les vitesses recueillies soient distribuées normalement avec une moyenne de 72 km/h et un écart-type de 8 km/h, quelle est approximativement la proportion d'automobiles ayant commis un excès de vitesse? **Exercice 6.7** Pour l'assemblage d'une machine, on prend

une longueur normale de moyenne 10 cm et d'écart-type 0,2 cm. On groupe les cylindres en 3 catégories: A: d'éfectueux et inutilisés; B: à 17 degrés de liberté; 3. d'une variable de Student à 8 degrés de liberté; 4. d'une variable de Fisher (uniquement

Tables statistiques

Tab.7.1-Table des quantiles d'une variable normale centrée réduite

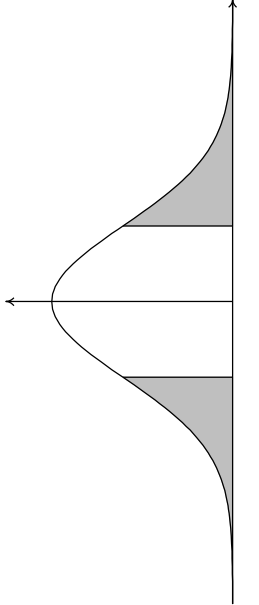


(Probabilité de trouver une valeur inférieure à)



7.9997.9998108

Tabele de valori ale funcției de densitate a variabilei aleatoare normale centrata și standardizate (u: valoarea variabilei, t: valoarea funcției de densitate)



$-u+u\alpha/2\alpha00.010.020.030.040.050.060.070.08$

$0-\infty+\infty$

0.31.03641.01520.99450.97410.95420.934

	0.09
0_{∞}	2.57582.32632.17012.05371.96001.880
0.1	1.64491.59821.55481.51411.47581.439
0.2	1.28161.25361.22651.20041.17501.150

Tab.7.4-Table des quantiles d'une variable χ^2 à n degrés de liberté

	ordre du quantile					
	0.01	0.025	0.05	0.95	0.975	0.99
$n=1$	0.000157	0.000982	0.003932	3.841	5.024	6.635
2	0.02010	0.05064	0.103	5.991	7.378	9.210
3	0.115	0.216	0.352	7.815	9.348	11.34
4	0.297	0.484	0.711	9.488	11.14	13.28
5	0.554	0.831	1.145	11.07	12.83	15.09
6	0.872	1.237	1.635	12.59	14.45	16.81
7	1.239	1.690	2.167	14.07	16.01	18.48
8	1.646	2.180	2.733	15.51	17.53	20.09
9	2.088	2.700	3.325	16.92	19.02	21.67
10	2.558	3.247	3.940	18.31	20.48	23.21
11	3.053	3.816	4.575	19.68	21.92	24.72
12	3.571	4.404	5.226	21.03	23.34	26.22
13	4.107	5.009	5.892	22.36	24.74	27.69
14	4.660	5.629	6.571	23.68	26.12	29.14
15	5.229	6.262	7.261	25.00	27.49	30.58
16	5.812	6.908	7.962		28.85	32.00
17	6.408	7.564	8.672		30.19	33.41
18	7.015	8.231	9.390		31.53	34.81
19	7.638	8.907	10.123		32.85	36.19
20	8.260	9.591	10.853			37.57
21	8.897	10.281	11.593			38.93
22	9.542	10.981	12.343			
23	10.201	11.691	13.093			
24	10.861	12.401	13.853			
25	11.521	13.121	14.613			
26	12.197	13.857	15.379			
27	12.878	14.598	16.153			
28	13.564	15.338	16.919			
29	14.253	16.077	17.689			
30	14.953	16.791	18.460			
31	15.653	17.501	19.232			
32	16.353	18.207	19.997			
33	17.053	18.910	20.758			
34	17.753	19.610	21.517			
35	18.453	20.307	22.275			
36	19.153	21.001	23.029			
37	19.853	21.692	23.781			
38	20.553	22.380	24.531			
39	21.253	23.065	25.279			
40	21.953	23.748	26.025			
41	22.653	24.428	26.769			
42	23.353	25.105	27.511			
43	24.053	25.779	28.251			
44	24.753	26.450	28.989			
45	25.453	27.118	29.725			
46	26.153	27.783	30.459			
47	26.853	28.445	31.191			
48	27.553	29.104	31.921			
49	28.253	29.760	32.649			
50	28.953	30.413	33.376			
51	29.653	31.063	34.101			
52	30.353	31.710	34.824			
53	31.053	32.354	35.545			
54	31.753	32.995	36.264			
55	32.453	33.633	36.981			
56	33.153	34.268	37.696			
57	33.853	34.900	38.409			
58	34.553	35.529	39.120			
59	35.253	36.155	39.829			
60	35.953	36.778	40.536			
61	36.653	37.398	41.241			
62	37.353	38.015	41.944			
63	38.053	38.629	42.645			
64	38.753	39.240	43.344			
65	39.453	39.848	44.041			
66	40.153	40.453	44.736			
67	40.853	41.055	45.429			
68	41.553	41.654	46.120			
69	42.253	42.250	46.809			
70	42.953	42.843	47.496			
71	43.653	43.433	48.181			
72	44.353	44.020	48.864			
73	45.053	44.604	49.545			
74	45.753	45.185	50.224			
75	46.453	45.763	50.901			
76	47.153	46.338	51.576			
77	47.853	46.910	52.249			
78	48.553	47.479	52.920			
79	49.253	48.045	53.589			
80	49.953	48.608	54.256			
81	50.653	49.168	54.921			
82	51.353	49.725	55.584			
83	52.053	50.279	56.245			
84	52.753	50.830	56.904			
85	53.453	51.378	57.561			
86	54.153	51.923	58.216			
87	54.853	52.465	58.869			
88	55.553	53.004	59.520			
89	56.253	53.540	60.169			
90	56.953	54.073	60.816			
91	57.653	54.603	61.461			
92	58.353	55.130	62.104			
93	59.053	55.654	62.745			
94	59.753	56.175	63.384			
95	60.453	56.693	64.021			
96	61.153	57.208	64.656			
97	61.853	57.720	65.289			
98	62.553	58.229	65.920			
99	63.253	58.735	66.549			
100	63.953	59.238	67.176			
101	64.653	59.738	67.801			
102	65.353	60.235	68.424			
103	66.053	60.729	69.045			
104	66.753	61.220	69.664			
105	67.453	61.708	70.281			
106	68.153	62.193	70.896			
107	68.853	62.675	71.509			
108	69.553	63.154	72.120			
109	70.253	63.630	72.729			
110	70.953	64.103	73.336			
111	71.653	64.573	73.941			
112	72.353	65.040	74.544			
113	73.053	65.504	75.145			
114	73.753	65.965	75.744			
115	74.453	66.423	76.341			
116	75.153	66.878	76.936			
117	75.853	67.330	77.529			
118	76.553	67.779	78.120			
119	77.253	68.225	78.709			
120	77.953	68.668	79.296			
121	78.653	69.108	79.881			
122	79.353	69.545	80.464			
123	80.053	69.978	81.045			
124	80.753	70.408	81.624			
125	81.453	70.835	82.201			
126	82.153	71.259	82.776			
127	82.853	71.679	83.349			
128	83.553	72.096	83.920			
129	84.253	72.510	84.489			
130	84.953	72.921	85.056			
131	85.653	73.329	85.621			
132	86.353	73.733	86.184			
133	87.053	74.134	86.745			
134	87.753	74.532	87.304			
135	88.453	74.927	87.861			
136	89.153	75.319	88.416			
137	89.853	75.708	88.969			
138	90.553	76.094	89.520			
139	91.253	76.477	90.069			
140	91.953	76.857	90.616			
141	92.653	77.234	91.161			
142	93.353	77.608	91.704			
143	94.053	77.979	92.245			
144	94.753	78.347	92.784			
145	95.453	78.712	93.321			
146	96.153	79.074	93.856			
147	96.853	79.433	94.389			
148	97.553	79.789	94.920			
149	98.253	80.142	95.449			
150	98.953	80.492	95.976			
151	99.653	80.839	96.501			
152	100.353	81.182	97.024			
153	101.053	81.522	97.545			
154	101.753	81.859	98.064			
155	102.453	82.192	98.581			
156	103.153	82.522	99.096			
157	103.853	82.849	99.609			
158	104.553	83.172	100.120			
159	105.253	83.492	100.629			
160	105.953	83.809	101.136			
161	106.653	84.122	101.641			
162	107.353	84.432	102.144			
163	108.053	84.739	102.645			
164	108.753	85.042	103.144			
165	109.453	85.342	103.641			
166	110.153	85.639	104.136			
167	110.853	85.933	104.629			
168	111.553	86.224	105.120			
169	112.253	86.512	105.609			
170	112.953	86.797	106.096			
171	113.653	87.079	106.581			
172	114.353	87.358	107.064			
173	115.053	87.634	107.545			
174	115.753	87.907	108.024			
175	116.453	88.177	108.501			
176	117.153	88.444	108.976			
177	117.853	88.708	109.449			
178	118.553	88.969	109.920			
179	119.253	89.227	110.389			
180	119.953	89.482	110.856			
181	120.653	89.734	111.321			
182	121.353	89.983	111.784			
183	122.053	90.229	112.245			
184	122.753	90.472	112.704			
185	123.453	90.712	113.161			
186	124.153	90.949	113.616			
187	124.853	91.183	114.069			
188	125.553	91.414	114.520			
189	126.253	91.642	114.969			
190	126.953	91.867	115.416			
191	127.653	92.089	115.861			
192	128.353	92.308	116.304			
193	129.053	92.524	116.745			
194	129.753	92.737	117.184			
195	130.453	92.947	117.621			
196	131.153	93.154	118.056			
197	131.853	93.358	118.489			
198	132.553	93.559	118.920			
199	133.253	93.757	119.349			
200	133.953	93.952	119.776			
201	134.653	94.144	120.201			

Tab.7.5-Table des quantiles d'une variable de Student`an de gésdeliberté

	ordre du quantile			
	0.95	0.975	0.99	0.995
$n=1$	6.314	12.71	31.82	63.66
2	2.920	4.303	6.965	9.925
3	2.353	3.182	4.541	5.841
4	2.132	2.776	3.747	4.604
5	2.015	2.571	3.365	4.032
6	1.943	2.447	3.143	3.707
7	1.895	2.365	2.998	3.499
8	1.860	2.306	2.896	3.355
9	1.833	2.262	2.821	3.250
10	1.812	2.228	2.764	3.169
11	1.796	2.201	2.718	3.106
12	1.782	2.179	2.681	3.055
13	1.771	2.160	2.650	3.012
14	1.761	2.145	2.624	2.977
15	1.753	2.131	2.602	2.947
16	1.746	2.120	2.583	2.921
17	1.740	2.110	2.567	2.898
18	1.734	2.101	2.552	2.878
19	1.729	2.093	2.539	2.861
20	1.725	2.086	2.528	2.845
21	1.721	2.080	2.518	2.831
22	1.718	2.075	2.509	2.819
23	1.715	2.071	2.501	2.808
24	1.712	2.067	2.494	2.798
25	1.710	2.064	2.487	2.789
26	1.708	2.061	2.481	2.781
27	1.706	2.058	2.475	2.774
28	1.704	2.056	2.469	2.767
29	1.702	2.054	2.464	2.761
30	1.701	2.052	2.459	2.756
31	1.700	2.050	2.454	2.751
32	1.699	2.048	2.449	2.746
33	1.698	2.046	2.444	2.741
34	1.697	2.044	2.439	2.736
35	1.696	2.042	2.434	2.731
36	1.695	2.040	2.429	2.726
37	1.694	2.038	2.424	2.721
38	1.693	2.036	2.419	2.716
39	1.692	2.034	2.414	2.711
40	1.691	2.032	2.409	2.706
41	1.690	2.030	2.404	2.701
42	1.689	2.028	2.399	2.696
43	1.688	2.026	2.394	2.691
44	1.687	2.024	2.389	2.686
45	1.686	2.022	2.384	2.681
46	1.685	2.020	2.379	2.676
47	1.684	2.018	2.374	2.671
48	1.683	2.016	2.369	2.666
49	1.682	2.014	2.364	2.661
50	1.681	2.012	2.359	2.656
51	1.680	2.010	2.354	2.651
52	1.679	2.008	2.349	2.646
53	1.678	2.006	2.344	2.641
54	1.677	2.004	2.339	2.636
55	1.676	2.002	2.334	2.631
56	1.675	2.000	2.329	2.626
57	1.674	1.998	2.324	2.621
58	1.673	1.996	2.319	2.616
59	1.672	1.994	2.314	2.611
60	1.671	1.992	2.309	2.606
61	1.670	1.990	2.304	2.601
62	1.669	1.988	2.299	2.596
63	1.668	1.986	2.294	2.591
64	1.667	1.984	2.289	2.586
65	1.666	1.982	2.284	2.581
66	1.665	1.980	2.279	2.576
67	1.664	1.978	2.274	2.571
68	1.663	1.976	2.269	2.566
69	1.662	1.974	2.264	2.561
70	1.661	1.972	2.259	2.556
71	1.660	1.970	2.254	2.551
72	1.659	1.968	2.249	2.546
73	1.658	1.966	2.244	2.541
74	1.657	1.964	2.239	2.536
75	1.656	1.962	2.234	2.531
76	1.655	1.960	2.229	2.526
77	1.654	1.958	2.224	2.521
78	1.653	1.956	2.219	2.516
79	1.652	1.954	2.214	2.511
80	1.651	1.952	2.209	2.506
81	1.650	1.950	2.204	2.501
82	1.649	1.948	2.199	2.496
83	1.648	1.946	2.194	2.491
84	1.647	1.944	2.189	2.486
85	1.646	1.942	2.184	2.481
86	1.645	1.940	2.179	2.476
87	1.644	1.938	2.174	2.471
88	1.643	1.936	2.169	2.466
89	1.642	1.934	2.164	2.461
90	1.641	1.932	2.159	2.456
91	1.640	1.930	2.154	2.451
92	1.639	1.928	2.149	2.446
93	1.638	1.926	2.144	2.441
94	1.637	1.924	2.139	2.436
95	1.636	1.922	2.134	2.431
96	1.635	1.920	2.129	2.426
97	1.634	1.918	2.124	2.421
98	1.633	1.916	2.119	2.416
99	1.632	1.914	2.114	2.411
100	1.631	1.912	2.109	2.406

Tab.7.6–Table des quantiles d'ordre 0.95 d'une variable de Fisher a_{n_1} et n_2 de degrés de liberté

	$n_1=12345$	6	7	8	9	10	12	14	16	20	30	∞
$n_2=11$	161.4199.5215.7224.6230.2234.0236.8238.9240.5241.9243.9245.4246.5248.0250.1254.3	218.5119.0019.1619.2519.3019.3319.3519.3719.3819.4019.4119.4219.4319.4519.4619.50	310.139.5529.2779.1179.0138.9418.8878.8458.8128.7868.7458.7158.6928.6608.6178.526	47.7096.9446.5916.3886.2566.1636.0946.0415.9995.9645.9125.8735.8445.8035.7465.628	56.6085.7865.4095.1925.0504.9504.8764.8184.7724.7354.6784.6364.6044.5584.4964.365	65.9875.1434.7574.5344.3874.2844.2074.1474.0994.0604.0003.9563.9223.8743.8083.669	75.5914.7374.3474.1203.9723.8663.7873.7263.6773.6373.5753.5293.4943.4453.3763.230	85.3184.4594.0663.8383.6873.5813.5003.4383.3883.3473.2843.2373.2023.1503.0792.928	95.1174.2563.8633.6333.4823.3743.2933.2303.1793.1373.0733.0252.9892.9362.8642.707			
	104.9654.1033.7083.4783.3263.2173.1353.0723.0202.9782.9132.8652.8282.7742.7002.538	114.8443.9823.5873.3573.2043.0953.0122.9482.8962.8542.7882.7392.7012.6462.5702.404	124.7473.8853.4903.2593.1062.9962.9132.8492.7962.7532.6872.6372.5992.5442.4662.296	134.6673.8063.4113.1793.0252.9152.8322.7672.7142.6712.6042.5542.5152.4592.3802.206	144.6003.7393.3443.1122.9582.8482.7642.6992.6462.6022.5342.4842.4452.3882.3082.131	154.5433.6823.2873.0562.9012.7902.7072.6412.5882.5442.4752.4242.3852.3282.2472.066	164.4943.6343.2393.0072.8522.7412.6572.5912.5382.4942.4252.3732.3332.2762.1942.010	174.4513.5923.1972.9652.8102.6992.6142.5482.4942.4502.3812.3292.2892.2302.1481.960	184.4143.5553.1602.9282.7732.6612.5772.5102.4562.4122.3422.2902.2502.1912.1071.917	194.3813.5223.1272.8952.7402.6282.5442.4772.4232.3782.3082.2562.2152.1552.0711.878		
	204.3513.4933.0982.8662.7112.5992.5142.4472.3932.3482.2782.2252.1842.1242.0391.843	214.3253.4673.0722.8402.6852.5732.4882.4202.3662.3212.2502.1972.1562.0962.0101.812	224.3013.4433.0492.8172.6612.5492.4642.3972.3422.2972.2262.1732.1312.0711.9841.783	234.2793.4223.0282.7962.6402.5282.4422.3752.3202.2752.2042.1502.1092.0481.9611.757	244.2603.4033.0092.7762.6212.5082.4232.3552.3002.2552.1832.1302.0882.0271.9391.733	254.2423.3852.9912.7592.6032.4902.4052.3372.2822.2362.1652.1112.0692.0071.9191.711	264.2253.3692.9752.7432.5872.4742.3882.3212.2652.2202.1482.0942.0521.9901.9011.691	274.2103.3542.9602.7282.5722.4592.3732.3052.2502.2042.1322.0782.0361.9741.8841.672	284.1963.3402.9472.7142.5582.4452.3592.2912.2362.1902.1182.0642.0211.9591.8691.654294.1833.3282.9342.7012.5452.	0772.0031.9481.9041.8391.7441.509504.0343.1832.7902.5572.4002.2862.1992.1302.0732.0261.9521.8951.8501.7841.6871		

Tab.7.7-Table des quantiles d'ordre 0.99 d'une variable de Fisher a_{n_1} et n_2 de degrés de liberté

	$n_1=12345$	6	7	8	9	10	12	14	16	20	30	∞
$n_2=1405250005403562557645859$	5928	5981	6022	6056	6106	6143	6170	6209	6261	6366		
	298.5099.0099.1799.2599.3099.3399.3699.3799.3999.4099.4299.4399.4499.4599.4799.50											
	334.1230.8229.4628.7128.2427.9127.6727.4927.3527.2327.0526.9226.8326.6926.5126.13											
	421.2018.0016.6915.9815.5215.2114.9814.8014.6614.5514.3714.2514.1514.0213.8413.46											
	516.2613.2712.0611.3910.9710.6710.4610.2910.1610.059.8889.7709.6809.5539.3799.020											
	613.7510.939.7809.1488.7468.4668.2608.1027.9767.8747.7187.6057.5197.3967.2296.880											
	712.259.5478.4517.8477.4607.1916.9936.8406.7196.6206.4696.3596.2756.1555.9925.650											
	811.268.6497.5917.0066.6326.3716.1786.0295.9115.8145.6675.5595.4775.3595.1984.859											
	910.568.0226.9926.4226.0575.8025.6135.4675.3515.2575.1115.0054.9244.8084.6494.311											
	1010.047.5596.5525.9945.6365.3865.2005.0574.9424.8494.7064.6014.5204.4054.2473.909											
	119.6467.2066.2175.6685.3165.0694.8864.7444.6324.5394.3974.2934.2134.0993.9413.602											
	129.3306.9275.9535.4125.0644.8214.6404.4994.3884.2964.1554.0523.9723.8583.7013.361											
	139.0746.7015.7395.2054.8624.6204.4414.3024.1914.1003.9603.8573.7783.6653.5073.165											
	148.8626.5155.5645.0354.6954.4564.2784.1404.0303.9393.8003.6983.6193.5053.3483.004											
	158.6836.3595.4174.8934.5564.3184.1424.0043.8953.8053.6663.5643.4853.3723.2142.868											
	168.5316.2265.2924.7734.4374.2024.0263.8903.7803.6913.5533.4513.3723.2593.1012.753											
	178.4006.1125.1854.6694.3364.1023.9273.7913.6823.5933.4553.3533.2753.1623.0032.653											
	188.2856.0135.0924.5794.2484.0153.8413.7053.5973.5083.3713.2693.1903.0772.9192.566											
	198.1855.9265.0104.5004.1713.9393.7653.6313.5233.4343.2973.1953.1163.0032.8442.489											
	208.0965.8494.9384.4314.1033.8713.6993.5643.4573.3683.2313.1303.0512.9382.7782.421											
	218.0175.7804.8744.3694.0423.8123.6403.5063.3983.3103.1733.0722.9932.8802.7202.360											
	227.9455.7194.8174.3133.9883.7583.5873.4533.3463.2583.1213.0192.9412.8272.6672.305											
	237.8815.6644.7654.2643.9393.7103.5393.4063.2993.2113.0742.9732.8942.7812.6202.256											
	247.8235.6144.7184.2183.8953.6673.4963.3633.2563.1683.0322.9302.8522.7382.5772.211											
	257.7705.5684.6754.1773.8553.6273.4573.3243.2173.1292.9932.8922.8132.6992.5382.169											
	267.7215.5264.6374.1403.8183.5913.4213.2883.1823.0942.9582.8572.7782.6642.5032.131											
	277.6775.4884.6014.1063.7853.5583.3883.2563.1493.0622.9262.8242.7462.6322.4702.097											
	287.6365.4534.5684.0743.7543.5283.3583.2263.1203.0322.8962.7952.7162.6022.4402.064297.5985.4204.5384.0453.7253.											
	8012.6652.5632.4842.3692.2031.805507.1715.0574.1993.7203.4083.1863.0202.8902.7852.6982.5622.4612.3822.2652.0981											

Listedestableaux

1.1 Codification de la variable Y	8
1.2 S'éri statistique de la variable Y	8
1.3 Tableau statistique complet....	9
3.1 Tableau des effectifs n_{jk}	44
3.2 Tableau des fréquences.....	45
3.3 Tableau des profils lignes.....	45
3.4 Tableau des profils colonnes.....	46
3.5 Tableau des effectifs théoriques n^*_{jk}	46
3.6 Tableau des écarts à l'indépendance e_{jk}	47
3.7 Tableau des e_{2jk}/n^*_{jk}	47
3.8 Tableau de contingence: effectifs n_{jk}	47
3.9 Tableau des fréquences f_{jk}	48
3.10 Tableau des profils lignes.....	48
3.11 Tableau des profils colonnes.....	48
3.12 Tableau des effectifs théoriques n^*_{jk}	48
3.13 Tableau des écarts à l'indépendance e_{jk}	48
3.14 Tableau des e_{2jk}/n^*_{jk}	48
3.15 Consommation de crème glacées.....	49
4.1 Tableau du prix d'un bien de consommation de 2000 à 2006.....	51
4.2 Tableau de l'indice simplifié du prix du tableau 4.1.....	51
4.3 Exemple: prix et quantité de trois biens pendant 3 ans.....	52
4.4 Mesures de l'inégalité entre pays.....	58
5.1 Biens manufacturés aux USA.....	60
5.2 Indices des prix à la consommation (France).....	62
5.3 Trafic d'un nombre de voyageurs SNCF.....	635
4D'ecomposition de la variable FRIG, méthode érietemporelle Prix moyen du Mazout pour 100 litres (achat entre 800 et 1500 litres) en CHF.....	806
'e.....	1107
5 Tables des quantiles d'une variable de Student à n degrés de liberté.....	1117
6 Tab.....	

7.7 Table des quantiles d'ordre 0.99 d'une variable de Fisher η_1 et η_2 de degrés de liberté	113
-----------------------------------------------------------------------------------------------------------------	-----

Table des figures

1.1 Diagramme en secteurs....	7
1.2 Diagramme en barres.....	8
1.3 Diagramme en secteurs des fréquences	9
1.4 Diagramme en barres des effectifs.	10
1.5 Diagramme en barres des effectifs cumulés	10
1.6 Diagramme en bâtonnets des effectifs pour une variable quantitative discrète	12
1.7 Fonction de répartition d'une variable quantitative discrète	13
1.8 Histogramme des effectifs.....	14
1.9 Histogramme des effectifs avec les deux dernières classes sautées.	15
1.10 Fonction de répartition d'une distribution groupée...	15
2.1 M'édiane quand n est impair.....	22
2.2 M'édiane quand n est pair.....	22
2.3 Asymétrie d'une distribution.....	28
2.4 Distributions mésokurtique et leptokurtique.....	28
2.5 Boîtes à moustaches pour la variable superficie en hectares (HApoly) des communes du canton de Neuchâtel.....	32
2.6 Boîtes à moustaches du "revenu moyen des habitants" des communes selon les provinces belges	33
3.1 L'usage de points.....	36
3.2 Exemples d'usage de points et coefficients de corrélation.....	38
3.3 L'usage de points, le résidu.....	38
3.4 L'adoption de la régression.....	40
4.1 Courbe de Lorenz.....	55
5.1 Dépenses en biens durables USA (milliards de dollars de 1982).....	61
5.2 Nombre de réfrigérateurs vendus de 1978 à 1985.....	61
5.3 Indices des prix à la consommation	65
5.4 Répartition de l'ordre 4 de la variable vente de réfrigérateurs'.....	68
5.5 Trafic d'un nombre de voyageurs SN	75
5.6 Evolution du prix du mazout en CHF (achat entre 800 et 1500), lissage exponentiel	76

6.3 Distribution d'une variable de Poisson avec $\lambda=1$.	94
6.4 Probabilité que la variable aléatoire soit inférieure à a .	95
6.5 Fonction de densité d'une variable uniforme.	97
6.6 Fonction de répartition d'une variable uniforme.	97
6.7 Fonction de densité d'une variable normale.	97
6.8 Fonction de répartition d'une variable normale.	98
6.9 Densité d'une normale centrée réduite, symétrie.	98
6.10 Fonction de densité d'une variable exponentielle avec $\lambda=1$.	100
6.11 Densité d'une variable de chi-carré avec $p=1, 2, \dots, 10$.	103
6.12 Densités de variables de Student avec $p=1, 2$ et 3 et d'une variable normale.	103
6.13 Densité d'une variable de Fisher.	104
6.14 Densité d'une normale bivariée.	104

Index

analyse combinatoire, 88
 arrangement, 89
 axiomatique, 84
 Bernoulli, 91
 bernoullienne, 91
 $\binom{n}{k}$, 91
 esaisonnalisation, 72
 diagramme en barres, 7
 de effectifs, 10
 n tonnets de effectifs, 12
 en boîte, 31
 en feuilles, 31

'esaisonnalisation,72diagrammeenbarres,7deseffectifs,10enb^atonnetsdeseffectifs,12enboite,31enfeuilles,

'esaisonnalisation,72diagrammeenbarres,7deseffectifs,10enb^atonnetsdeseffectifs,12enboite,31enfeuilles,

'epartition,12,15,21discontinue,23118

identification,52,99
'erieup,27
forwardoperator,66
fr'equence,6groupe,29histogrammedeseffectifs,14 pond'ee,21,30

op'rateur
 avance,66
 ded'ecalage,66
 dediff'ERENCE,66
 forward,66
 identit'e,66
 lag,66
 retard,66

param'etres
 d'aplatissement,28
 dedispersion,24
 deforme,27
 deposition,17
 marginaux,36

percentile,23

permutation
 avec'r'ep'etition,89
 sans'r'ep'etition,88

piechart,7probabilit'e,83,84conditionnelleetind'

 ependance,87
 th'eor`emedesprobabilit'es,87
 profilscolonnes,45lignes,45propri'et'es,102propri'et'esdes

quantile,23,36,106,107,109-111quartile,23quintile,23sharerat

des profils linéaires, 175 statistique, 7, 11, 12 tendance, 64 linéaire, 64, 66 logistique, 64 parabolique, 64 polynomiale, 64
 système de coordonnées, 84 ordinales, 5, 6 ordinales, 5, 7, 9, 11, 12, 13, 14, 15, 35 continue, 5, 12 discrète, 5, 11 uniforme, 5, 11
 tableau de contingence, 44
 de fréquences, 44
 des profils colonnes, 45
 de régression, 41, 42
 marginale, 36, 100
 propriétés, 101
 résiduelle, 42, 43